# Intégration et interrogation de données
## Askomics

F. Legeai

INRA

27 septembre 2018

# Intégration et interrogation de données
## Askomics

F. Legeai

INRA

27 septembre 2018

# Outline

# Big data and the deluge of life science data

### Big data

Datasets so **large** or **complex** that traditional data processing is inadequate [Laney2001]

### Life science: data deluge [Aldhous1993]

- computerized biomedical data
- Genomics and bioinformatics

Science. 1993 Oct 22;262(5133):502-3.

**Managing the genome data deluge.**

Aldhous P.

PMID: 8211171 [PubMed - indexed for MEDLINE]

Science. 1995 Aug 4;269(5224):630.

**Europe opens institute to deal with gene data deluge.**

Williams N.

PMID: 7624788 [PubMed - indexed for MEDLINE]

# Complexity of life science data: distributed

- 1500+ biological databases [Galperin2015]
- Lack of interoperability
- Some efforts of unified access (BioMart, InterMine...)

# What to expect for 2025?

Our estimation is that genomics is one of the most demanding domain in terms of

- data acquisition
- data storage
- data distribution
- data analysis

## Big Data: Astronomical or Genomical?

Zachary D. Stephens[1], Skylar Y. Lee[1], Faraz Faghri[2], Roy H. Campbell[2], Chengxiang Zhai[3], Miles J. Efron[4], Ravishankar Iyer[1], Michael C. Schatz[5]*, Saurabh Sinha[3]*, Gene E. Robinson[6]*
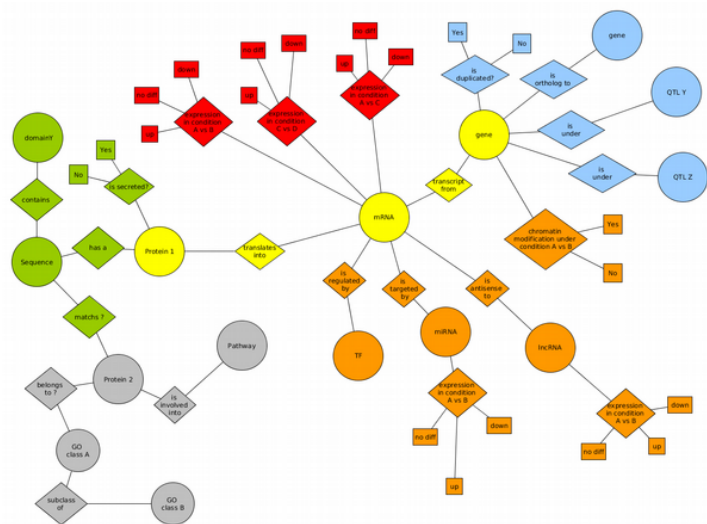
PLOS Biology | DOI:10.1371/journal.pbio.1002195   July 7, 2015

**Table 1. Four domains of Big Data in 2025.** In each of the four domains, the projected annual storage and computing needs are presented across the data lifecycle.

| Data Phase | Astronomy | Twitter | YouTube | Genomics |
|---|---|---|---|---|
| **Acquisition** | 25 zetta-bytes/year | 0.5–15 billion tweets/year | 500–900 million hours/year | 1 zetta-bases/year |
| **Storage** | 1 EB/year | 1–17 PB/year | 1–2 EB/year | 2–40 EB/year |
| **Analysis** | In situ data reduction | Topic and sentiment mining | Limited requirements | Heterogeneous data and analysis |
| | Real-time processing | Metadata analysis | | Variant calling, ~2 trillion central processing unit (CPU) hours |
| | Massive volumes | | | All-pairs genome alignments, ~10,000 trillion CPU hours |
| **Distribution** | Dedicated lines from antennae to server (600 TB/s) | Small units of distribution | Major component of modern user's bandwidth (10 MB/s) | Many small (10 MB/s) and fewer massive (10 TB/s) data movement |

doi:10.1371/journal.pbio.1002195.t001

# Data analysis needs data integration

# Data everywhere!

# Data are distributed

### Definition: Entity

Anything that can be identified (i.e. the things we can talk about)

- some informations about an entity in a repository
- other informations about the same entity in another repository

your question requires to combine entity descriptions from multiple datasets

Only possible if:

- Entities are identified
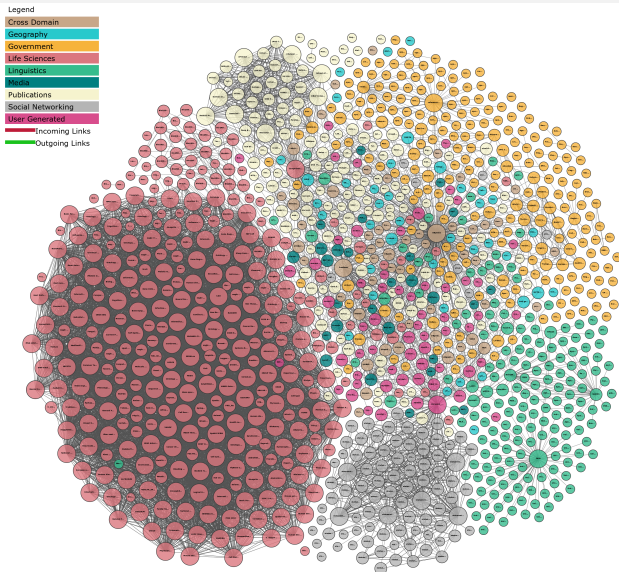- Datasets use the same identifier to describe the same entity

Good luck with your spreadsheets!  :-)

# Data description involves hierarchies

- Knowledge has been formalized (e.g. in ontologies)
- Query engines (more or less) gracefully handle
  - linking data and ontologies
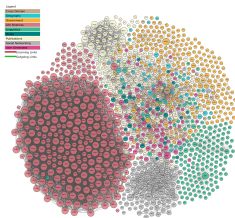  - reasoning on the ontologies

## The good news (1/3)

Semantic Web technologies (RDF, SPARQL, OWL) address all of these

# Linked open data (in 2017)



Linked open data cloud, by M. Schmachtenberg, C. Bizer, A. Jentzsch and R. Cyganiak http://lod-cloud.net/

# Linked Open Data



The not so good news: Linked data are here... but still have to be adopted by end users

"Real" users

- do not contribute (yet) their data to the LOD cloud
- do not use the LOD cloud for analyzing their own data (yet)

"Real" users need to retain their ability to interact with **their** data

# The RDF format

- Identify things
    - Describe them
    - their characteristics
    - their relations to other things
    - their categories
- Possibly combine descriptions from multiple places
- Elements of the descriptions should themselves be identiable
- Support rich querying and reasoning on descriptions



"Now! *That* should clear up a few things around here!"

# Resources (aka entities)

## Resource (= entity)

Resource: anything in the universe

## Two kinds of resource

- Referent: resource with an identity
  - identified by an **URI or IRI)**
  - there can be 0, 1 or more IRI identifying the same resource
- Litteral: resource defined by its value
  - represented as a string (e.g. "foo" or 'foo')
  - each litteral has 1 **datatype** that defines its nature and the range of possible values
    - default=string
    - string, integer, float, date, dateTime, boolean...
    - use XSD

# Litteral

- Represented by a string (delimited by single or double-quotes; escape if necessary)
- Its **datatype** (default = string) explains how to interpret
  - "12" = the string with two characters ('1' and '2')
  - "12"∧∧xsd:integer = the string representing the number twelve
  - "12.0"∧∧xsd:float = the string representing the number twelve
  - "true" = the string with four characters
  - "true"∧∧xsd:boolean = the string representing the boolean true
  - "1"∧∧xsd:boolean = the string representing the boolean true
  - "2017-08-24T18:54:42"∧∧xsd:dateTime
- String litterals can be qualified by a **language tag**
  - "ADN"@fr or "DNA"@en

# RDF triple

RDF triple = <subject, predicate, object>

- **subject**: the resource being described
- **predicate**: the relation (from the subject to the object)
- **object**: the value of the predicate for the subject
  (one of them if the predicate can have several values for the subject;
  in this case use as many triples as necessary)

# RDF triple

- an IRI can be the subject of many triples
- an IRI can the the object of many triples
- an IRI can be the subject of some triples and the object of other triples
- a litteral can only be an object

# RDF graph

RDF graph = set of RDF triples

- Set
  - No order among the triples
  - Redundant triples are ignored
- Structure = labeled directed graph
  - nodes = resources (subject or object of as least one triple)
  - edges = triples
    - start = subject IRI
    - end = object IRI
    - edge label = property IRI

The graph itself can be seen as a resource (**named graph**)

- identified by its IRI
- useful for describing who generated the graph, when, how, from which data...
- not mandatory (default graph)

# RDF in practice

RDF can be stored:

- in files (several formats: n3, turtle, xml-rdf,...)
- in triplestores ($\sim$ DB): Virtuoso, fuseki, sesame & RDF4J, Oracle, 4store,...
- within HTML: RDFa

RDF is supported by most languages (and most triplestores have bindings)

# RDF format: N-triples

- 1 triple / line
  - URI delimited by angle brackets ("<" and ">")
  - subject, predicate and object separated by space
  - line ends with a colon (".")
- no order between triples
- order within triple matters (obviously)

Straightforward, but not easy to read/write (for humans)

# RDF format: N-triples

```
1
2   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#firstName> "Frank" .
3   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#lastName> "Herbert" .
4   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#citizenOf>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#USA> .
5   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#authorOf>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Dune> .
6   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH>
•   <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Author> .
7   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH>
•   <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Person> .
8
9   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Dune>
•   <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
•   <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Book> .
10  <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Dune>
```

# RDF format: Turtle

- N-triples syntax still possible
- comments start with "#" until end of line
- prefix declaration for CURIE
  - @PREFIX prefixName:  <URI template> .
  - (also a triple)
  - usually at begining of file
- 1 triple / line
  - URIs can be full URIs (<...>) or CURIE (no delimiters)
  - a triple can mix both
  - line ends with ";" ⇒ next line has
    - same subject (no need to repeat)
  - line ends with "," ⇒ next line has
    - same subject (no need to repeat)
    - same property (no need to repeat)

# RDF format: Turtle

```
1   @prefix sf: <http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#> .
2   @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
3
4   #
5   # FRANK HERBERT
6   #
7   sf:FH sf:firstName "Frank" .
8   sf:FH sf:lastName "Herbert" .
9   sf:FH sf:citizenOf sf:USA .
10  sf:FH sf:authorOf sf:Dune .
11  sf:FH rdf:type sf:Author .
12  sf:FH rdf:type sf:Person .
13
14  sf:Dune rdf:type sf:Book .
15  sf:Dune rdf:type sf:SciFiBook .
16  sf:Dune sf:title "Dune" .
17  sf:Dune sf:wonAward sf:Nebula1965 .
18  sf:Dune sf:wonAward sf:Hugo1966 .
19
20  sf:Nebula1965 rdf:type sf:NebulaAward .
21  sf:Hugo1966 rdf:type sf:HugoAward .
```

# Other formats

### N3

Do not confuse with N-Triples!

- N3 = turtle + additional features
    - Named graphs
    - Inference rules
    - ...
- not a W3C recommendation

### RDF-XML

- XML serialization of RDF
- a major pain to read or edit these by hand

# RDF format: RDF-XML

```xml
1   <?xml version="1.0" encoding="utf-8"?>
2   <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
·     xmlns:sf="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#">
3     <rdf:Description rdf:about="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH">
4       <sf:firstName>Frank</sf:firstName>
5     </rdf:Description>
6     <rdf:Description rdf:about="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH">
7       <sf:lastName>Herbert</sf:lastName>
8     </rdf:Description>
9     <rdf:Description rdf:about="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH">
10      <sf:citizenOf rdf:resource="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#USA"/>
11    </rdf:Description>
12    <rdf:Description rdf:about="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH">
13      <sf:authorOf rdf:resource="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Dune"/>
14    </rdf:Description>
15    <rdf:Description rdf:about="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH">
16      <rdf:type rdf:resource="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Author"/>
17    </rdf:Description>
18    <rdf:Description rdf:about="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#FH">
19      <rdf:type rdf:resource="http://www.irisa.fr/dyliss/public/odameron/SciFiAwards#Person"/>
20    </rdf:Description>
```

# SPARQL Endpoints

- SPARQL queries are executed against RDF datasets, consisting of RDF graphs.
- A SPARQL endpoint accepts queries and returns results via HTTP.
- Endpoints can be generic (i.e.web accessible datasets) or specific to a dataset
- The results of SPARQL queries can be returned and/or rendered in a variety of formats (XML, JSON, RDF, HTML)

# SPARQL queries

SPARQL = query language similar to SQL
Variable names start with a question mark

```
1  @prefix foaf: <http://xmlns.com/foaf/0.1/> .
2  @prefix mus: <http://www.irisa.fr/dyliss/public/mus#> .
3
4  SELECT ?player ?instr
5   WHERE {
6     ?player foaf:name Alice
7     ?player mus:plays ?instr .
8   }
```

# SPARQL queries (2)

```
1  @prefix foaf: <http://xmlns.com/foaf/0.1/> .
2  @prefix mus: <http://www.irisa.fr/dyliss/public/mus#> .
3
4  SELECT ?player ?instr
5   WHERE {
6     ?player foaf:name Alice
7     ?player mus:plays ?instr
8     ?instr mus:instrclass corde .
9   }
```

# SPARQL queries (3)

```
SELECT ?player ?instr
 WHERE {
    ?player foaf:name Alice
    ?player mus:plays ?instr
    ?instr mus:instrclass corde
    ?player info:age ?age
    FILTER (?age <= 24) .
     }
```

# SPARQL queries (4)

+ many other features
  - OPTIONAL
  - UNION
  - ORDER BY
  - DISTINCT
  - . . .

# AskOmics: reconciling domain experts with Semantic Web

AskOmics is usefull for:

- Integrating data
- Querying data

# AskOmics (1) integrating data with RDF

## Integrating data

- Import your data files
    - CSV or TSV
    - RDF
    - GFF
- Import public knowledge bases (GO, Reactome, NCBI taxon,...)
- Declare (remote) SPARQL endpoints (in progress)

## The good news (2/3)

Don't worry about RDF, AskOmics generates it from your csv

# AskOmics (2) generating user-friendly SPARQL queries

### Querying data

- Graph-based user-friendly SPARQL query composer
  - based on an abstraction of your data
    - depends on the data structure (small)
    - not on the data themselves (possibly huge)
  - can be enriched (shortcuts and virtual links)
  - modular design (select/deselect datasets)
- Span multiple SPARQL endpoints (in progress)
- You do not have to see the SPARQL code
- Save the query result (obviously) in RDF or TSV

### The good news (3/3)

Don't worry about SPARQL queries, AskOmics generates them for you
(and runs them as well)

# Querying data

# AskOmics: new features in latest release

- Improved distribution
  - docker-compose
  - Genostack (thank you, GenOuest)
- Collaborative mode with public/private graphs management (authentification)
- Internal Modules
  - Gene Ontology
  - BioPax-v3 (Reactome,...)
- External modules integration
  - Metabolic Pathways (generated by Aureme)
  - Expression Analysis (generated by Askor)

# Managing your own data

# AskoR - The scheme



Inputs/ Outputs

# AskoR - 3 cons

Contrastes, contextes et conditions

# AskoR - Inputs

## Inputs / Outputs

- Comptages

| Geneid | T1A_1 | T1A_2 | T1A_3 | T1A_4 | T1K_1 |
|--------|-------|-------|-------|-------|-------|
| LOC100569617 | 380 | 12933 | 13406 | 3136 | 8825 |
| LOC100160354 | 2679 | 15 | 13 | 5 | 7 |
| **LOC100162386** | 21215 | 870 | 458 | 207 | **360** |

- Fichier de description des échantillons

| ID | condition | stage | treatment |
|----|-----------|-------|-----------|
| ID | condition | stage | treatment |
| T1A_1 | T1A | T1 | Acetone |
| T1A_2 | T1A | T1 | Acetone |
| T1A_3 | T1A | T1 | Acetone |
| T1A_4 | T1A | T1 | Acetone |
| T1K_1 | T1K | T1 | Kinoprene |
| T1K_2 | T1K | T1 | Kinoprene |
| T1K_3 | T1K | T1 | Kinoprene |
| T1K_4 | T1K | T1 | Kinoprene |
| T2A_1 | T2A | T2 | Acetone |
| T2A_2 | T2A | T2 | Acetone |
| T2A_3 | T2A | T2 | Acetone |
| T2A_4 | T2A | T2 | Acetone |
| T2K_1 | T2K | T2 | Kinoprene |
| T2K_2 | T2K | T2 | Kinoprene |
| T2K_3 | T2K | T2 | Kinoprene |
| T2K_4 | T2K | T2 | Kinoprene |

- Contrastes

| Condition | T1AvsT1K | T2AvsT2K | T1vsT2 | AvsK |
|-----------|----------|----------|--------|------|
| Condition | T1AvsT1K | T2AvsT2K | T1vsT2 | AvsK |
| T1A | + | | 0 | + | + |
| T1K | - | | 0 | + | - |
| T2A | | 0 | + | - | + |
| T2K | | 0 | - | - | - |

# AskoR - Commands

AskoR

Ligne de commande

Galaxy



https://github.com/askomics/askoR

# AskoR - Graphics

# AskoR - Results

# AskoR - Integration in Askomics (1)

# AskoR - Integration in Askomics (2)

# Askomics now

## AskOmics now

- Within Dyliss: data integration
- Dyliss ecosystem:
  - INRA
  - INSERM
  - Sanofi
  - received interest from: Fealinx + Theranexus
- Growing adoption by a wider user base
  - Docker image
  - GenoStack (support from GenOuest)
  - interface with Galaxy (developped by INRA)
  - AskoR (developped by INRA)
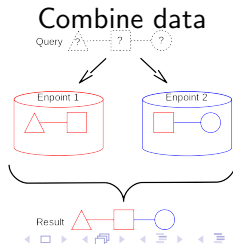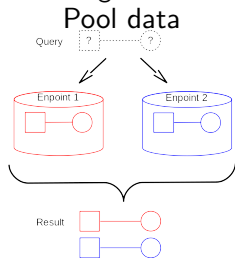
# The future of AskOmics

The future of AskOmics

- extend query expressivity
- improve the GUI ergonomy
- improve support for multiple endpoints
  - constellation of AskOmics
  - benefit from the FederatedQueryScaler PRE (with Wimmics)
- reach a larger user base

# AskOmics received contributions from:

- Meziane Aite
- Arnaud Belcour (intern)
- Charles Bettembourg (postdoc + Sanofi)
- Anthony Bretaudeau (INRA)
- Yvanne Chaussin (intern)
- Olivier Dameron (Univ. Rennes 1)
- Aurélie Evrard (postdoc INRA)
- Olivier Filangi (INRA)
- Xavier Garnier (intern + INRA)
- Maël Kerbiriou (intern)
- Fabrice Legeai (INRA)
- Sylvaine Bitteur (INRA / Agrocampus Ouest) for the logo
- Colleagues from Nantes for their insight on federated queries
  - Pascal Molli
  - Patricia Serrano Alvarado
  - Hala Skaf

# Federated queries principle

- Linked data
  - RDF repositories can be queried in SPARQL via endpoints
  - data from one endpoint can make references to data from another endpoint
- Federated queries span several endpoints
  - the SPARQL engine is responsible for propagating the query and merging the results
  - good news: supported by SPARQL language + query engines
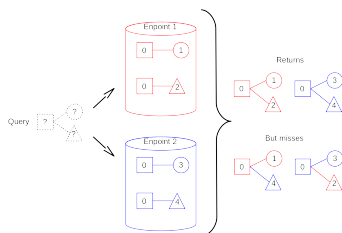  - not so good news: performances :-(

# Federated queries difficulty: endpoints not independent



Treating the endpoints independently fails when combining data
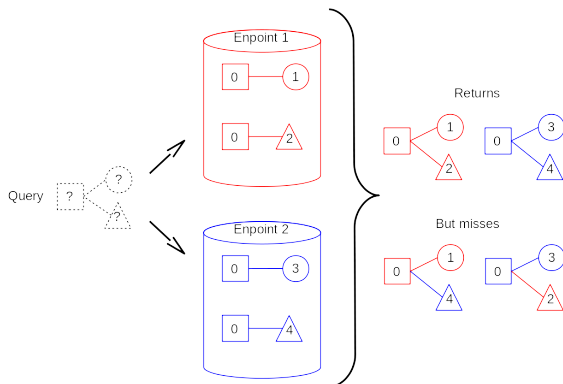
# Federated queries difficulty: endpoints can not be merged



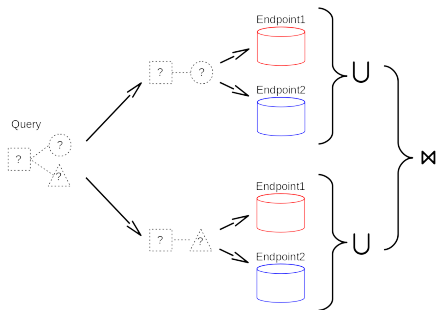Merging the endpoints is not a viable solution either

- each endpoint is potentially big
- merging
    - increases network traffic
    - increases storage consumption
    - decreases query answering performances
    - does not scale up to LOD

# Federated queries difficulty: queries must be decomposed



The query must be decomposed into fragments

so that the endpoints are not treated as if they were disjoint

# Federated queries: q. fragmentation increases complexity



Sending each fragment to each endpoint results in

- many subqueries for each endpoint (distant server overload)

- many unions and joins (local engine overload)

- potential transfer of large quantities of data before performing the joins, even if it ultimately few results (network overload)

# Processing federated queries: general approach

## Decompose the query into fragments

The fewer fragments the better: reduces joins

## For each fragment, select the relevant endpoints

The fewer endpoints the better (but no false negatives!): reduces joins

## Determine the order for processing the fragments (q. planning)

Start by the most selectives, maybe parallelize, and potentially rewrite the subqueries

These three aspects can be inter-dependent

# Processing federated queries: existing solutions

- General low level SPARQL engines: Virtuoso, Fuseki,...
    - well adapted for serving 1 endpoint
    - naive approach for federated queries
- Dedicated high level engines (FedX, HiBISCuS,...)
    - front end to low-level engines
    - hot domain (in general)
    - still not up to Life Science queries