

# 关联数据系统开发实例与 平台详解

陈 涛

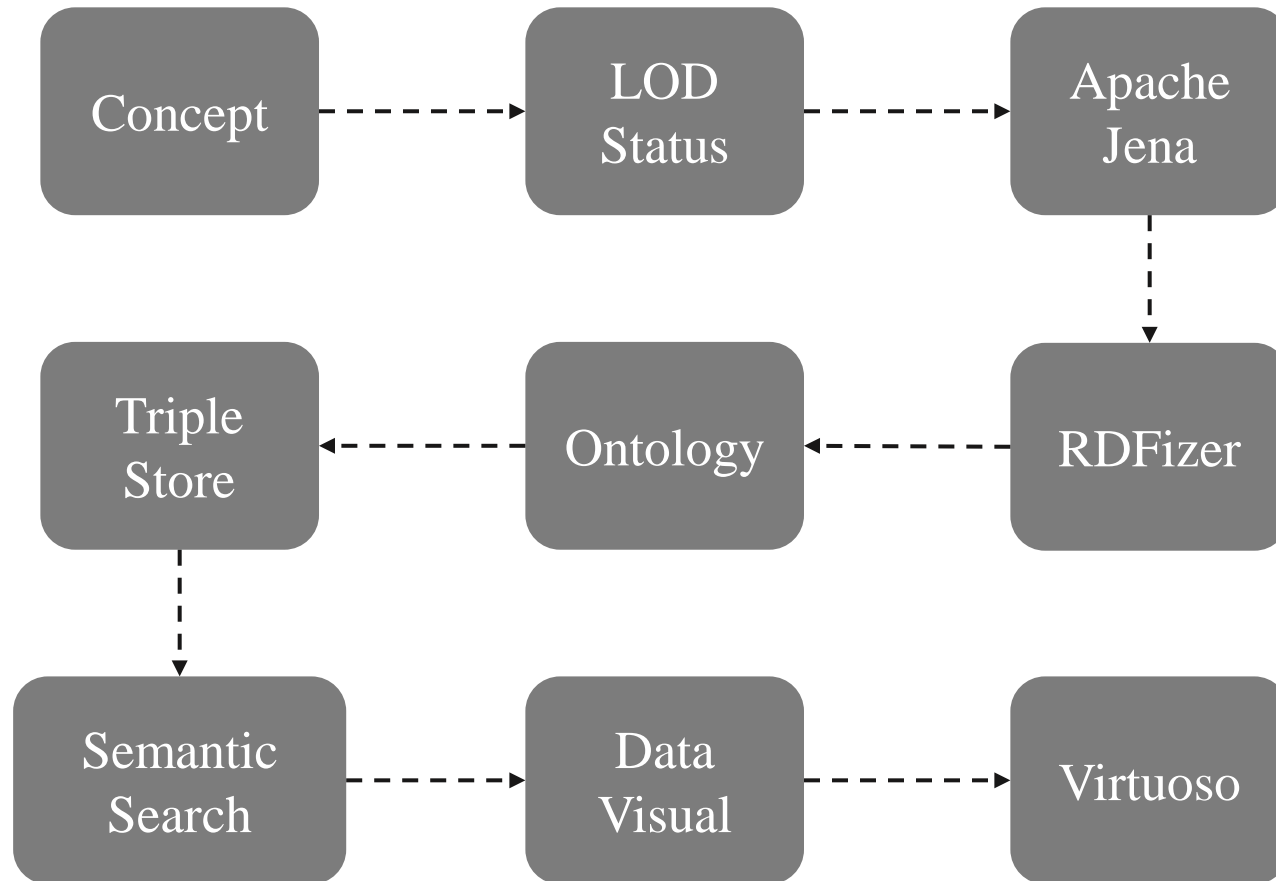
中国科学院上海生命科学信息中心

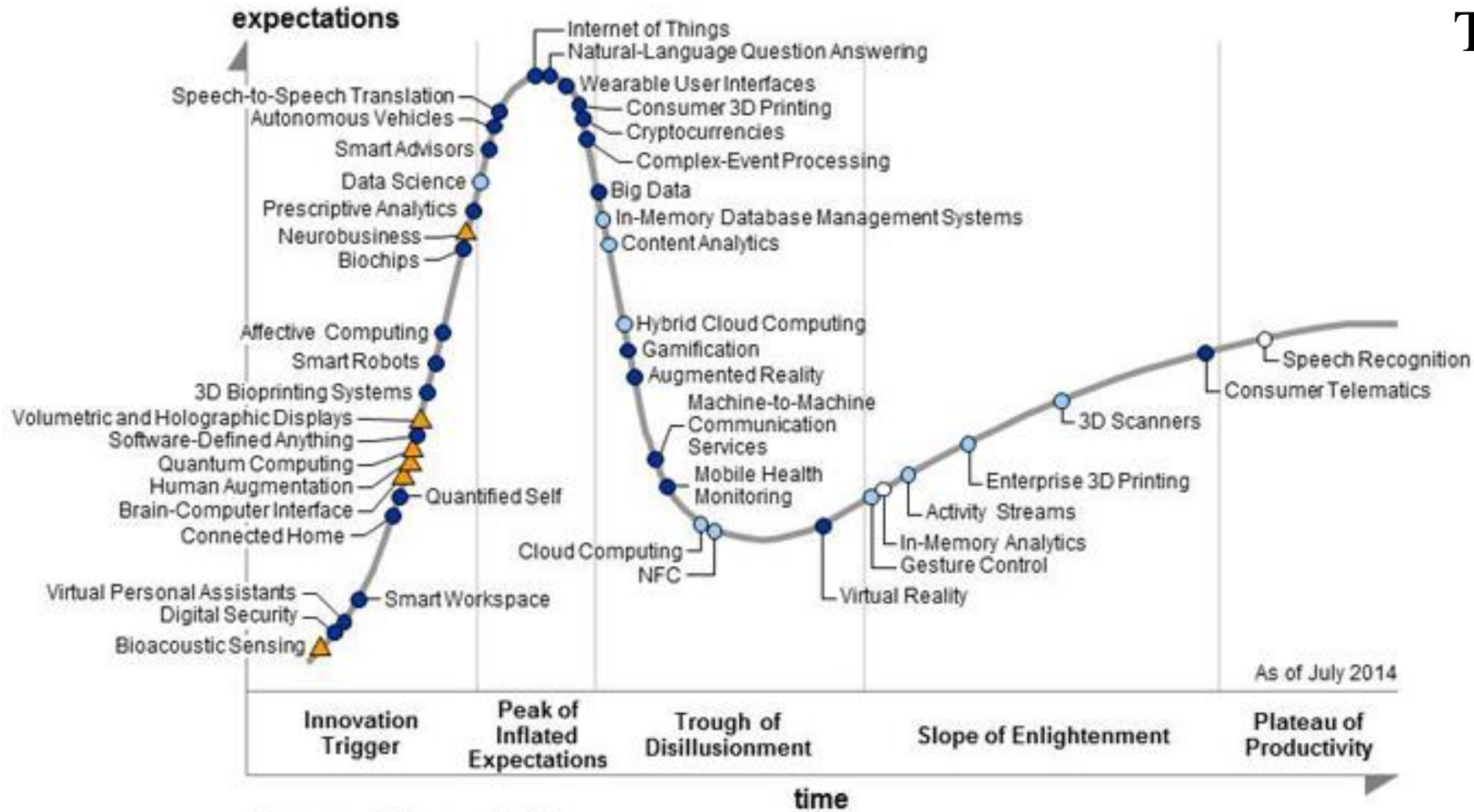
2015-05-05

QQ群：三人行（语义有你） [150461365](https://www.qq.com/group/150461365)



# Agenda





Plateau will be reached in:

- less than 2 years
- ◉ 2 to 5 years
- 5 to 10 years
- ▲ more than 10 years
- ⊗ obsolete before plateau

# The Internet of

# Things

# Everything



# Valuation of IoE



It's estimated that by the year 2020 more than fifty billion “things” will be connected to the internet.

**50,000,000,000**

How much is the Internet of everything worth?  
Cisco says \$19 trillion.

**19,000,000,000,000**

manufacturing (15%) ,healthcare (15%) ,insurance (11%)



# Semantic Web

- Semantic web is a **Web of data**.
- The term “Semantic Web” refers to W3C’s vision of the Web of linked data.
- It is about common formats for **integration** and **combination** of data drawn from diverse sources.
- It is also about language for recording how the data **relates** to real world objects. That allows a person, or a machine, to **start off in one database**, and then **move through an unending set of databases** which are connected not by wires but by being **about the same thing**.



# Linked Data



WIKIPEDIA  
The Free Encyclopedia

A term used to describe a recommended best practice for exposing, sharing, and connecting pieces of *data*, *information*, and *knowledge* on the Semantic Web using URIs and RDF.

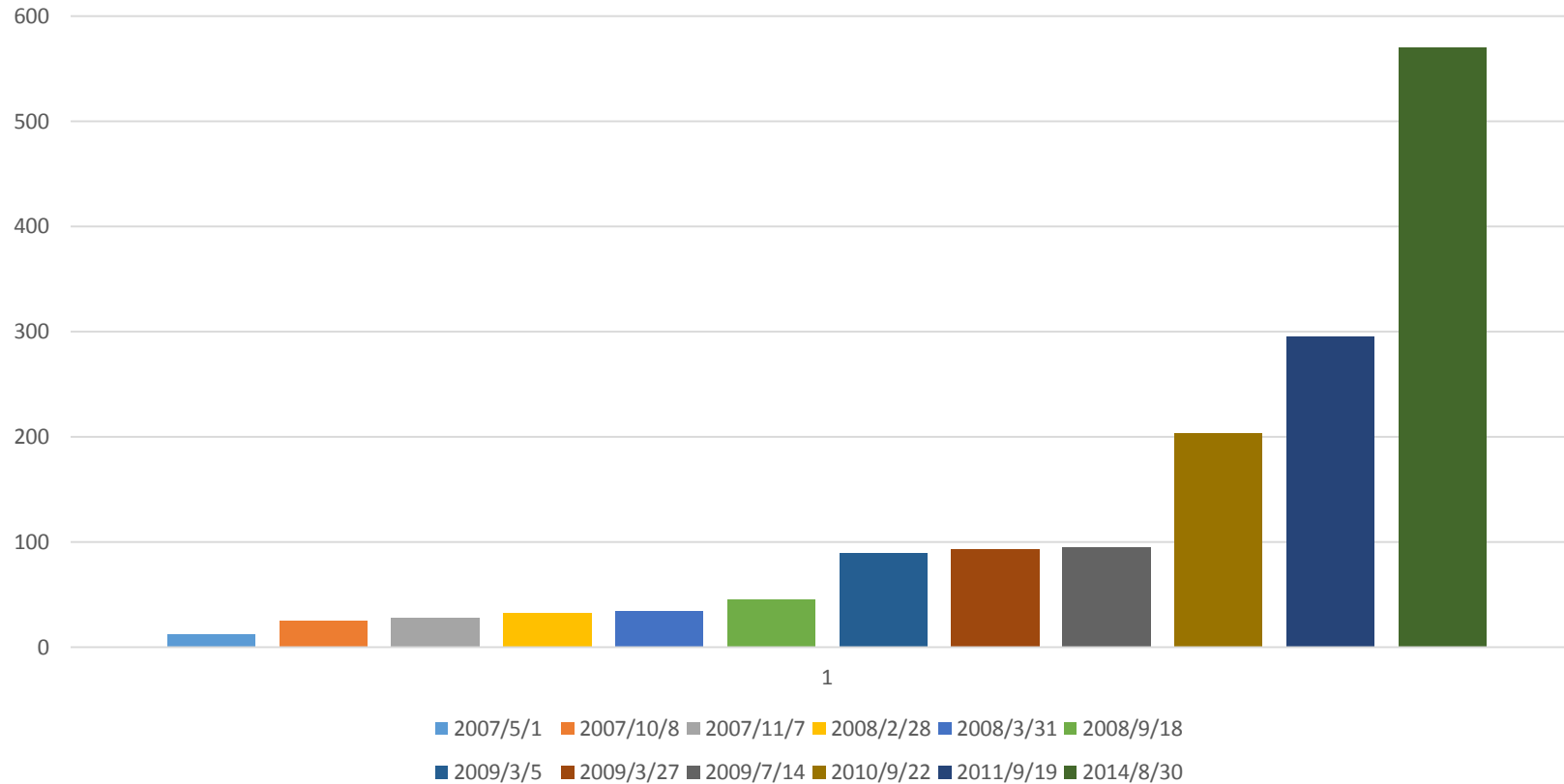
Connect Distributed Data across the Web





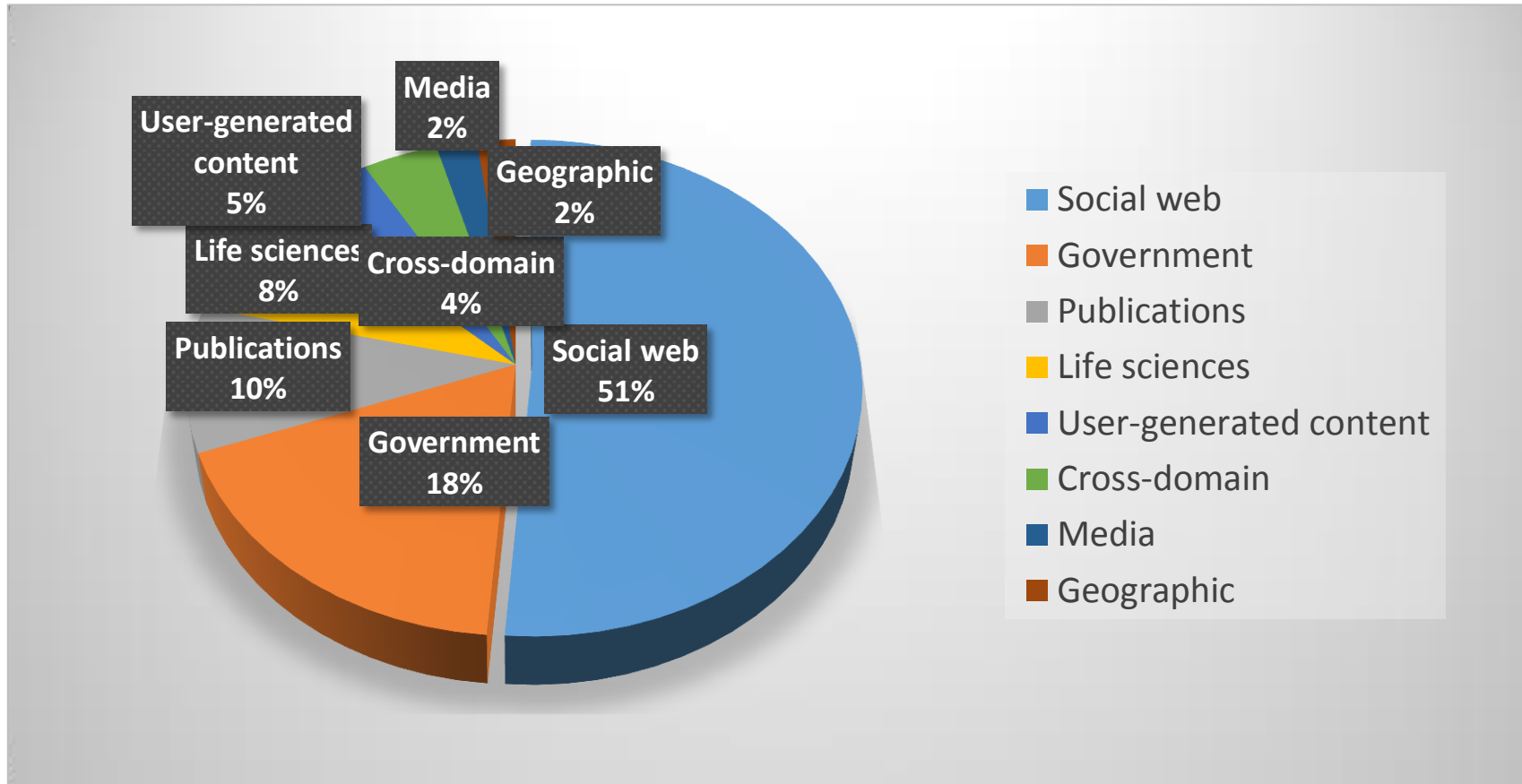
# Linked Open Data Cloud

## LOD Datasets





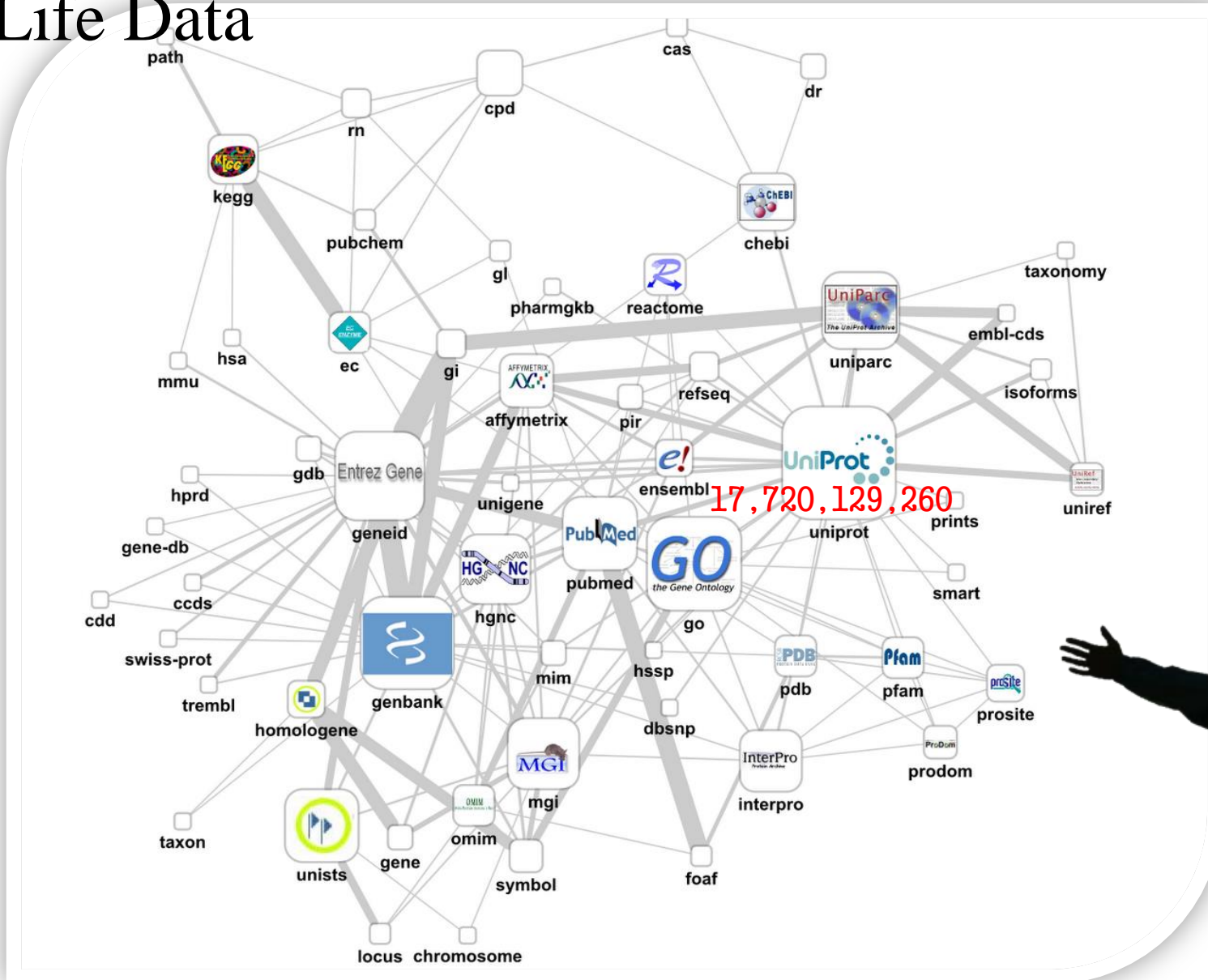
# LOD Cloud Status







# Linked Life Data





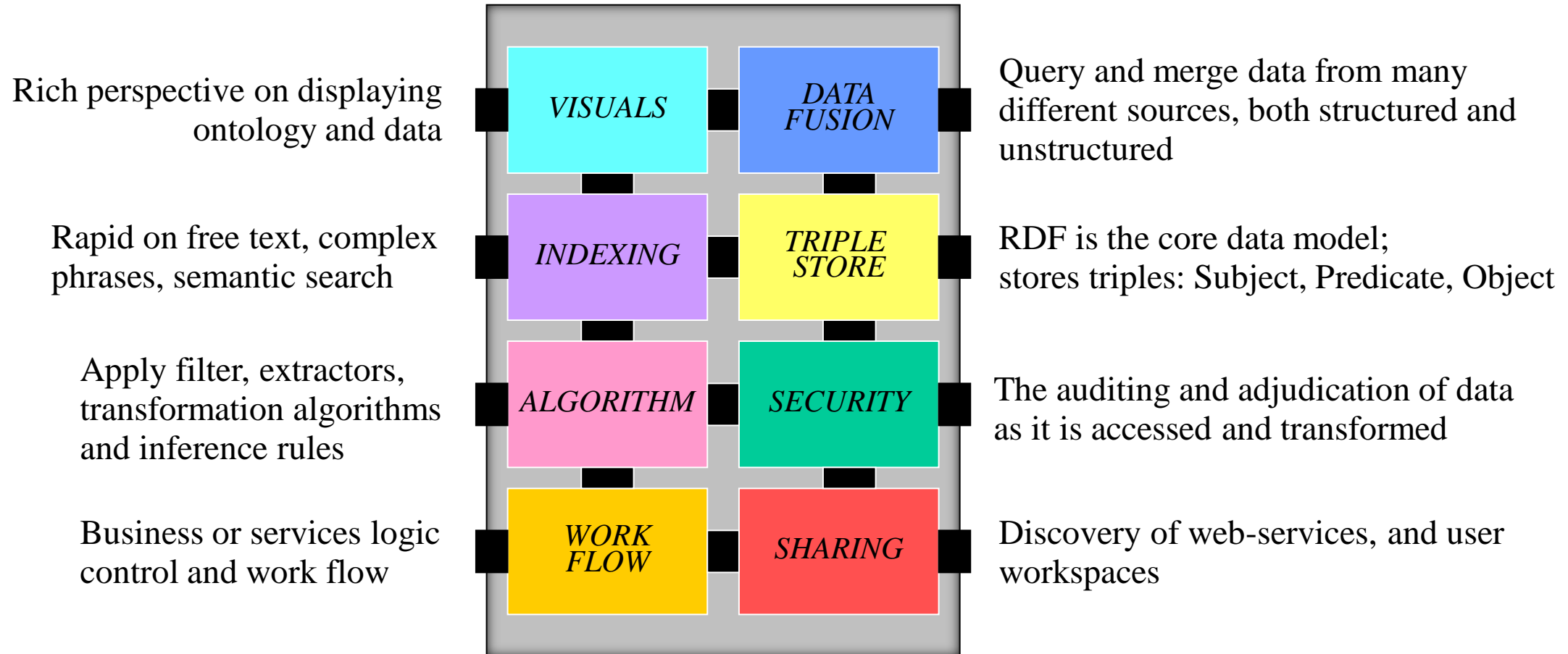
# Use Cases for SW & LD







# Core Components





## Framework for *SW* & *LD*



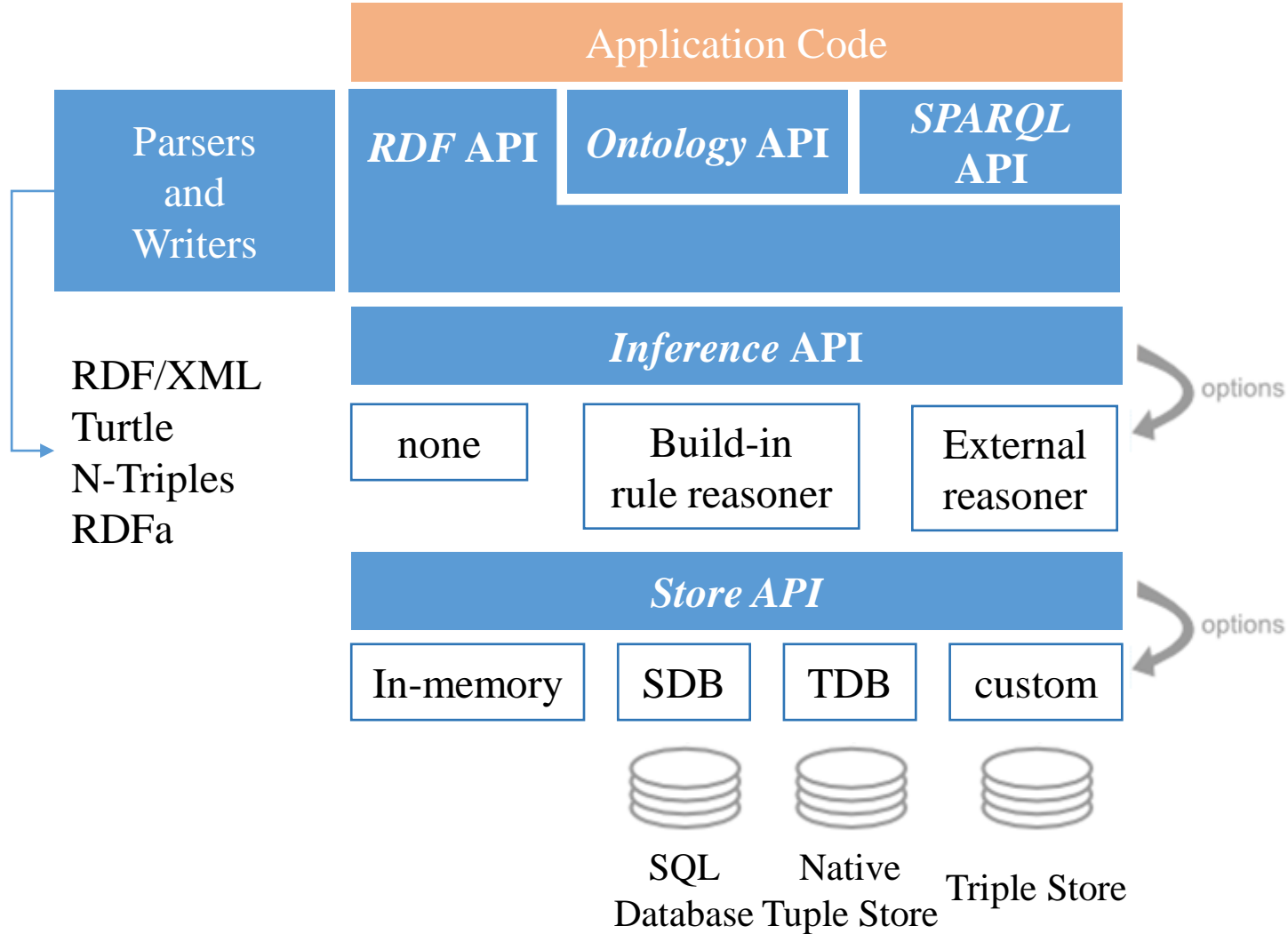
**Apache Jena** (or Jena in short) is a free and open source Java framework for building semantic web and Linked Data applications.



**OpenRDF Sesame** is a powerful Java framework for processing and handling RDF data. This includes creating, parsing, storing, inferencing and querying over such data.

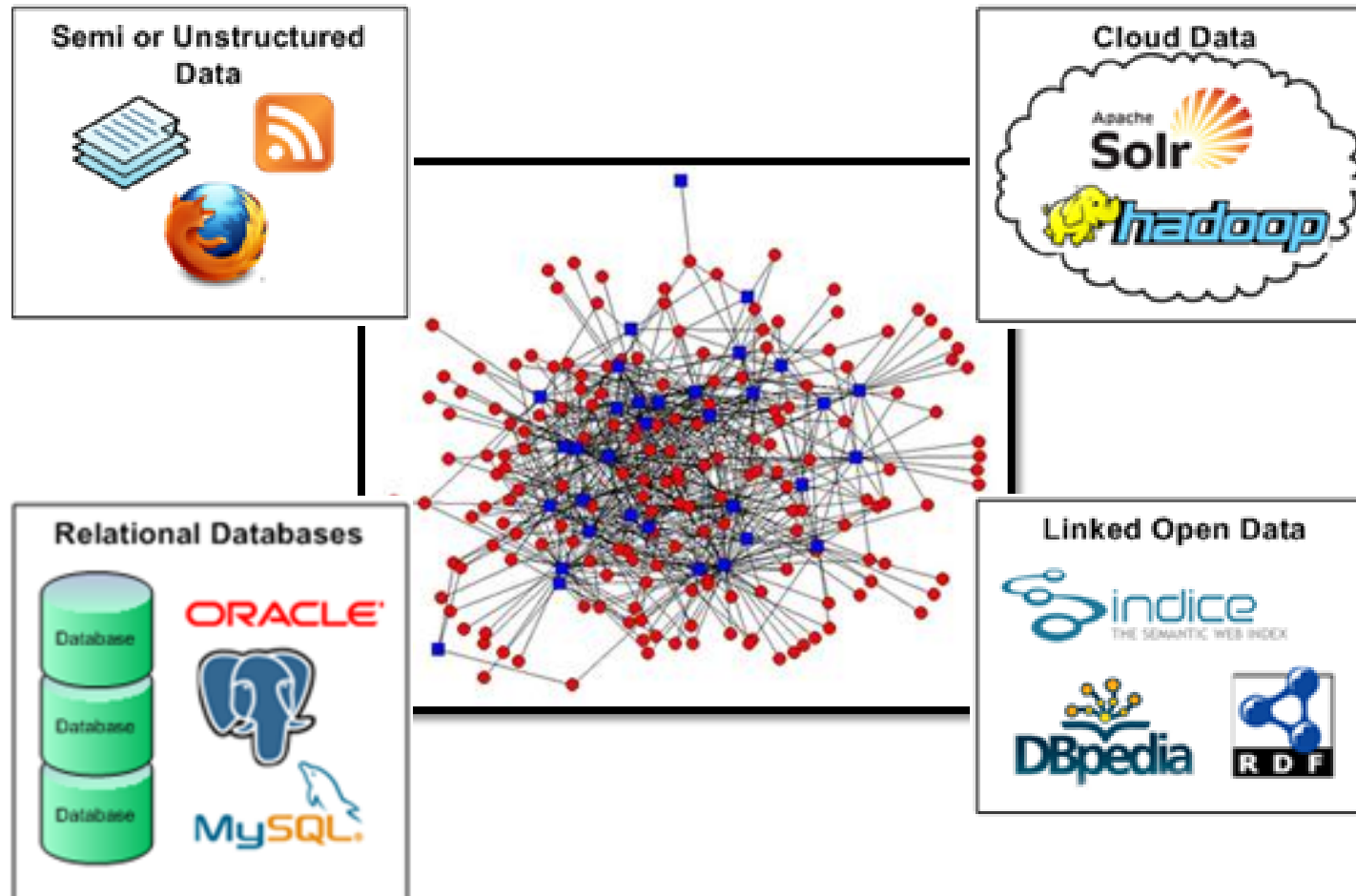


# Apache Jena





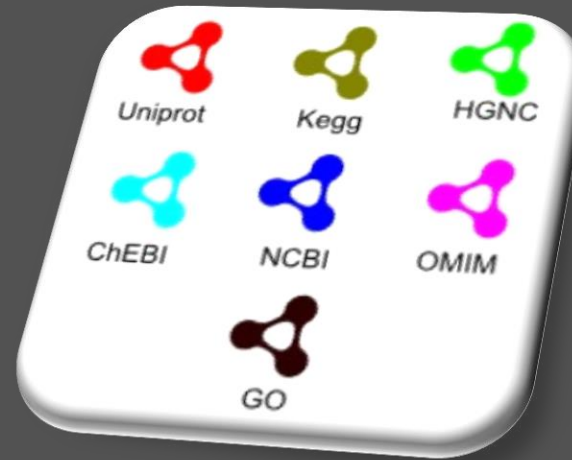
# Data Fusion / Integration Points





# Bio2RDF Process

- Dataset: 35 Triples: 11,900,533,153 Entities: 1,108,204,952



RDFization

Triple Store

URI Normalization

SERVICES: <http://hgnc.bio2rdf.org/sparql> <http://omim.bio2rdf.org/sparql>



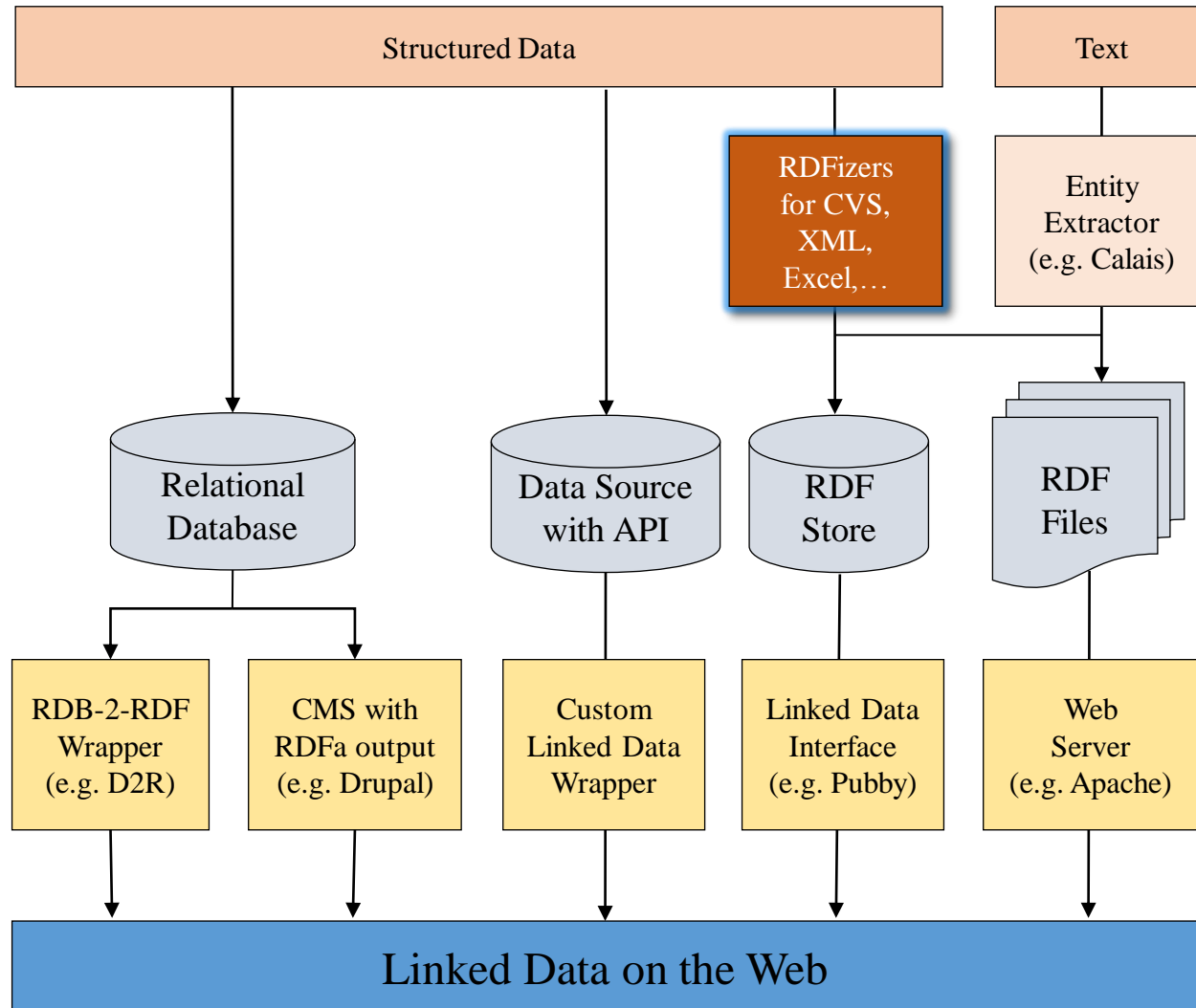


# Fusion Process

- 1** Download the original data
- 2** Transform the data into RDF
- 3** Keeping a link to the original data in the new RDF document
- 4** Store the converted database in a triple store
- 5** Give services for these using a RESTful interface



# Transformation / RDFization



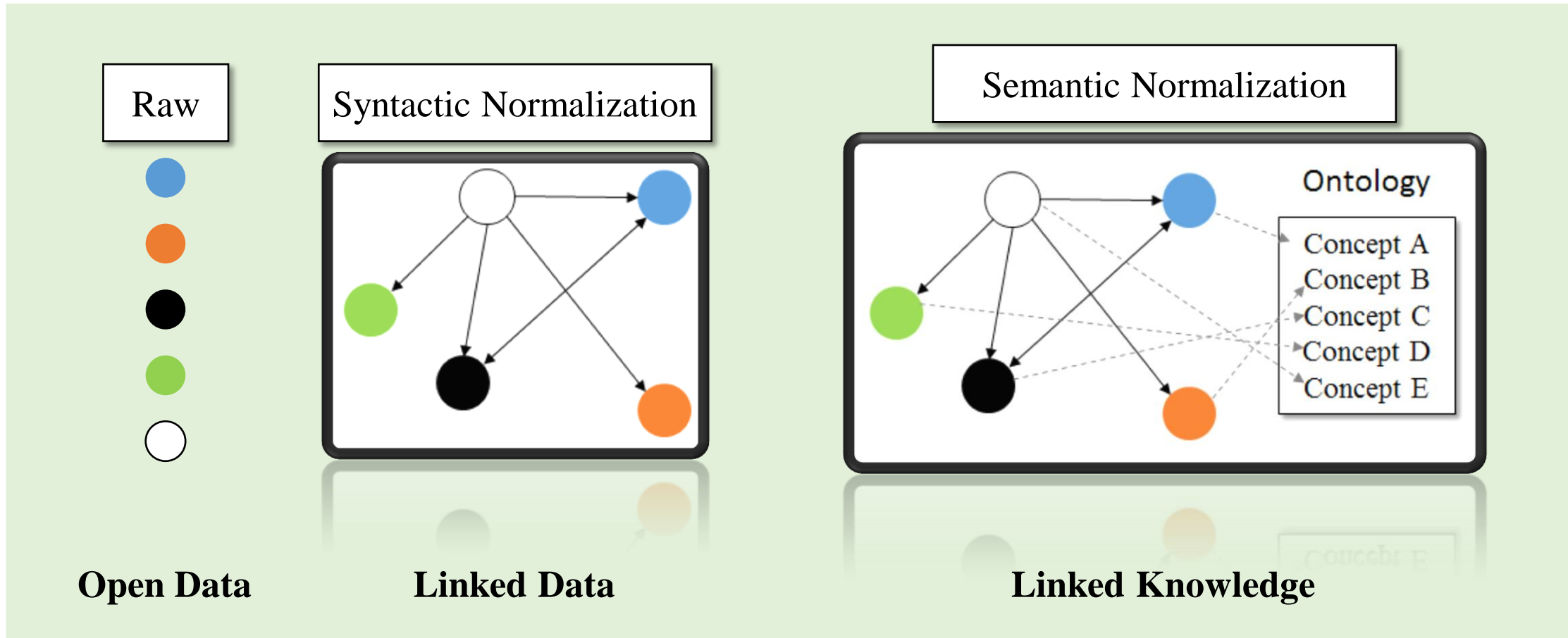


# Convert to RDF Tools

Format	Tools
<b>BibTex</b>	<ul style="list-style-type: none"><li>• BibBase</li><li>• bibtex2rdf</li><li>• Java BibTex-To-RDF Converter</li></ul>
<b>SQL</b>	<ul style="list-style-type: none"><li>• D2RQ</li><li>• OpenLink Virtuoso</li><li>• R2RML</li></ul>
<b>Excel</b>	<ul style="list-style-type: none"><li>• XLWrap</li><li>• Excel2rdf</li><li>• Sheet2RDF</li></ul>
<b>MARC</b>	<ul style="list-style-type: none"><li>• easyM2R</li><li>• marcmods2rdf</li></ul>
<b>XML</b>	<ul style="list-style-type: none"><li>• GRDDL</li><li>• SPARQL2XQuery</li></ul>
<b>Microformats</b>	<ul style="list-style-type: none"><li>• Any23</li></ul>

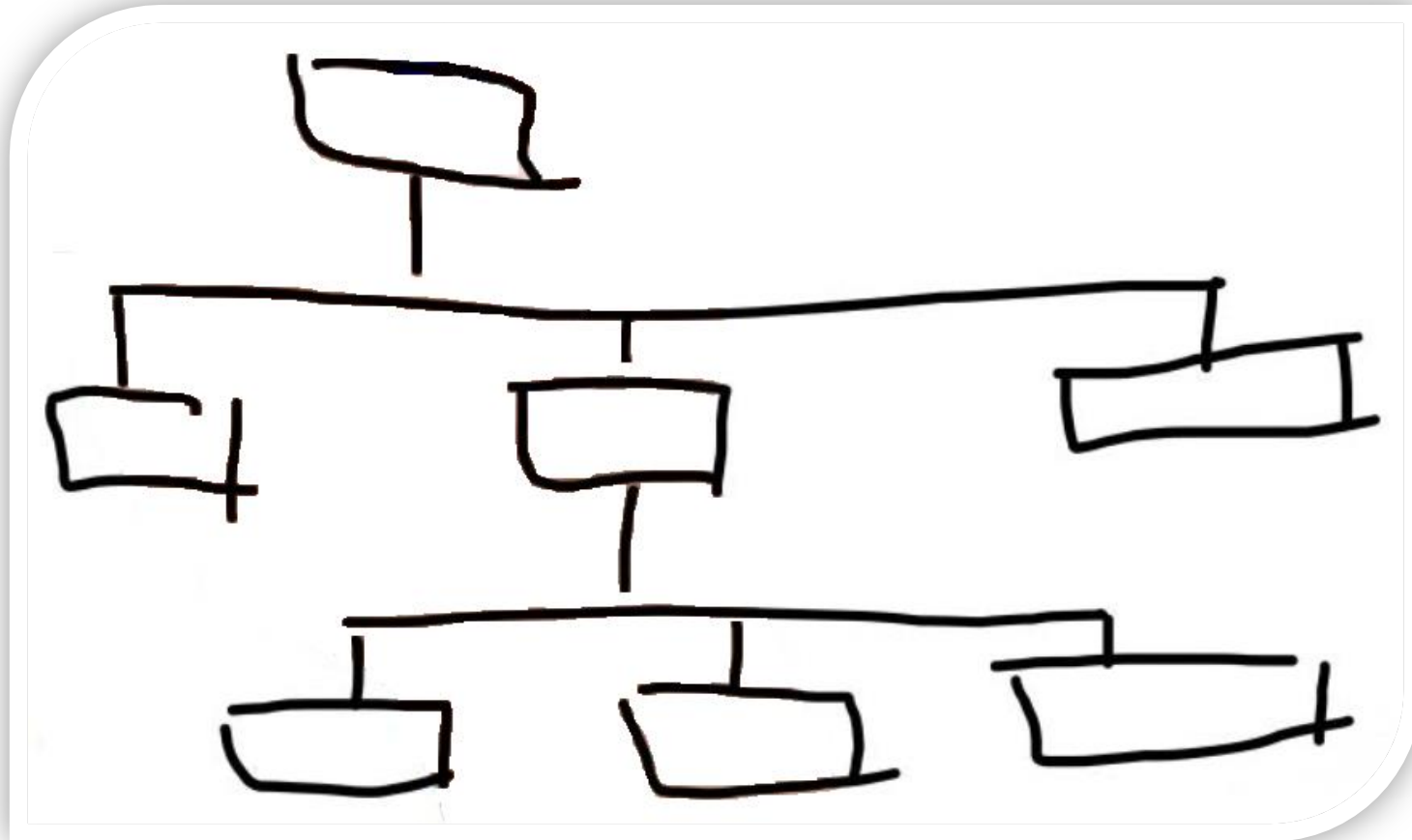


# RDF is only the First step, not End





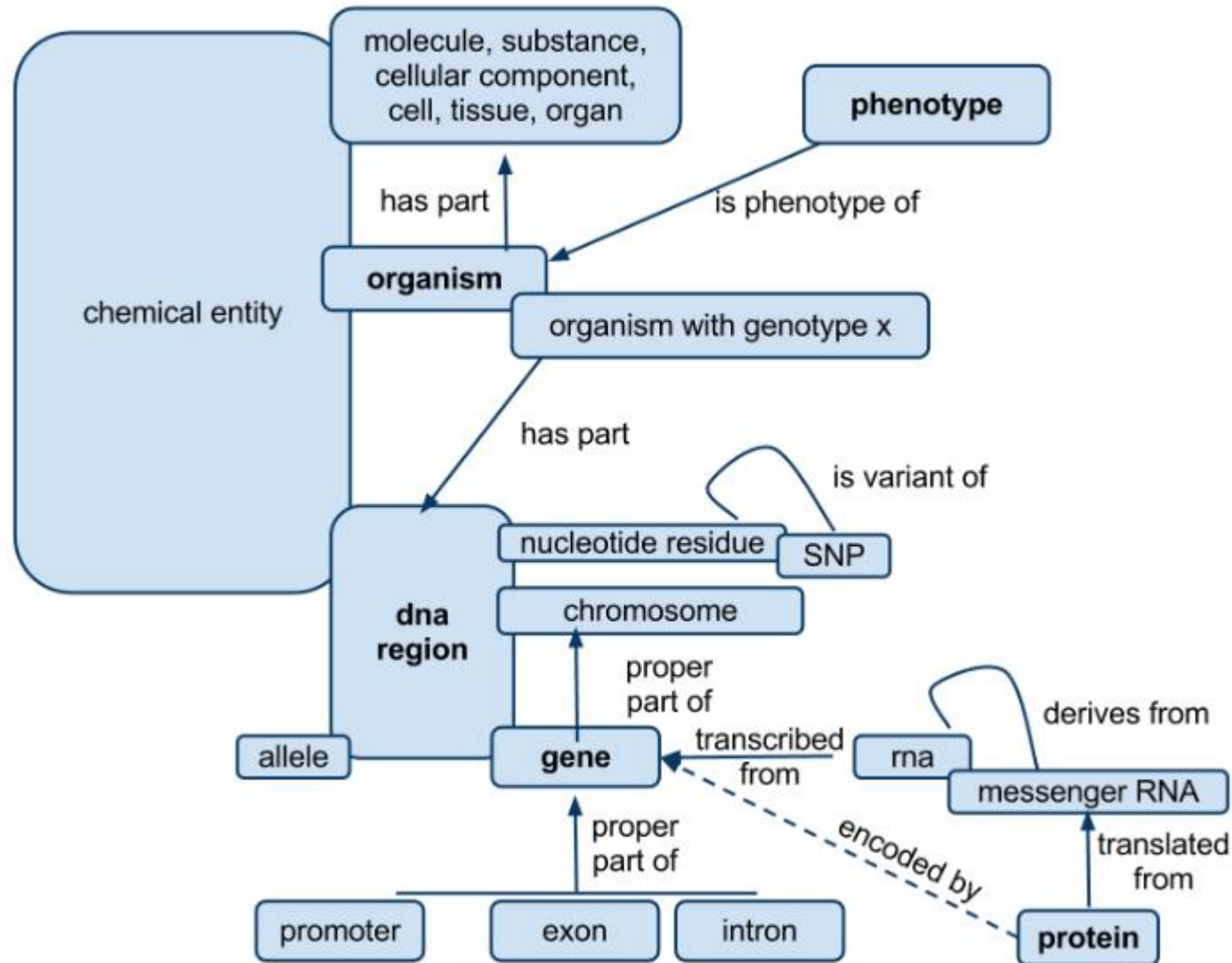
# Ontology



**ontology as a  
strategy to formally  
represent and  
integrate knowledge**

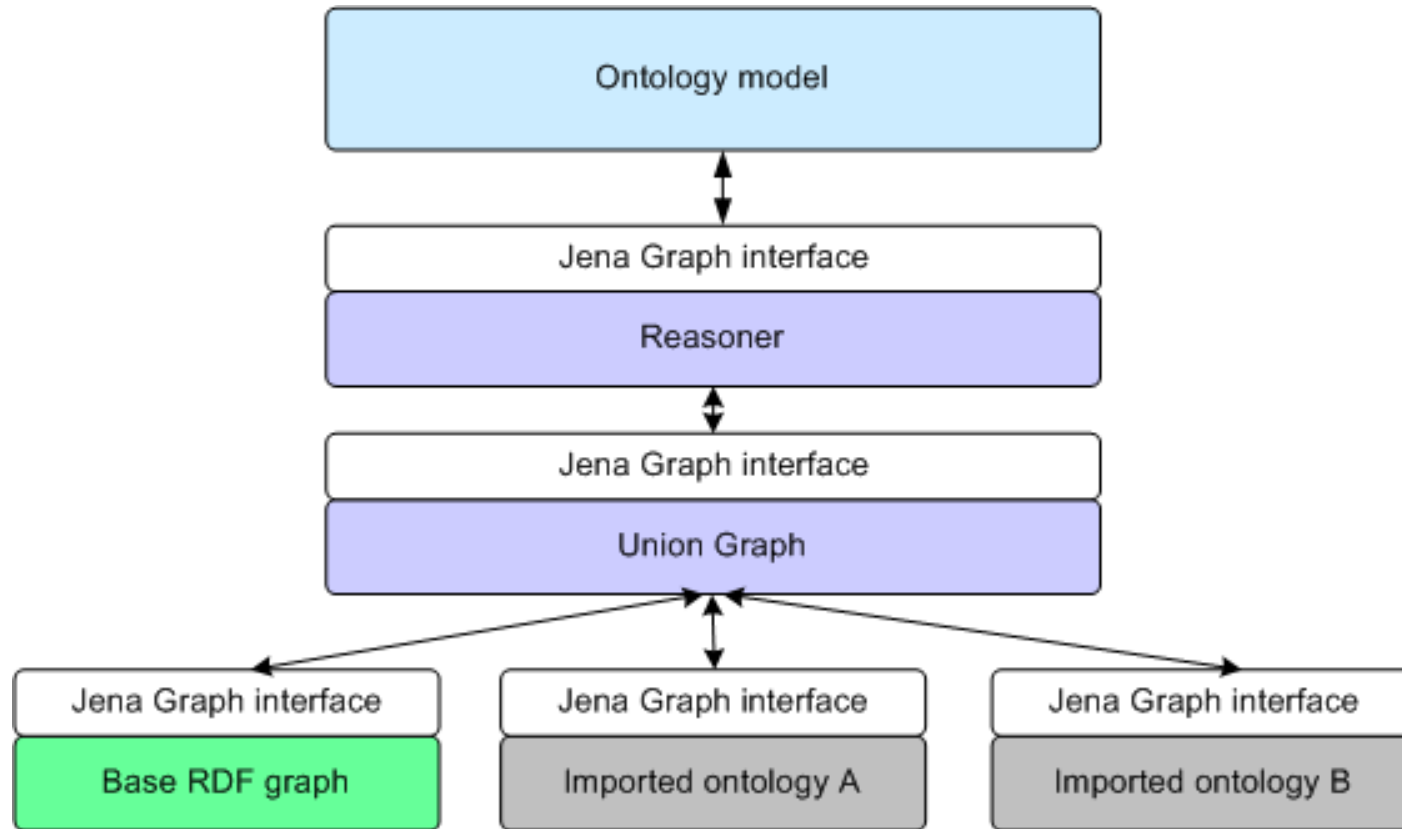


# Ontology for the representations of diverse biomedical knowledge





# Jena Ontology API



- Jena OntModel
- RDF Graph - SPARQL

Ontology model compound document structure for imports



# Jena Inference API

- **Transitive reasoner:** Provides support for storing and traversing class and property lattices. This implements just the *transitive* and *reflexive* properties of *rdfs:subPropertyOf* and *rdfs:subClassOf*.
- **RDFs rule reasoner:** Implements a configurable subset of the RDFs entailments.
- **OWL, OWL Mini, OWL Micro Reasoners:** A set of useful but incomplete implementation of the OWL/Lite subset of the OWL/Full language.
- **Generic rule reasoner:** A rule based reasoner that supports user defined rules. Forward chaining, tabled backward chaining and hybrid execution strategies are supported.



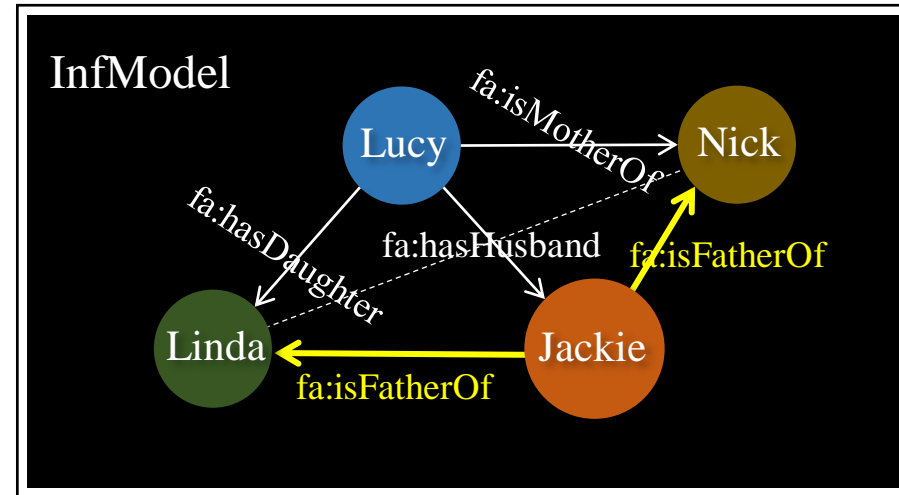
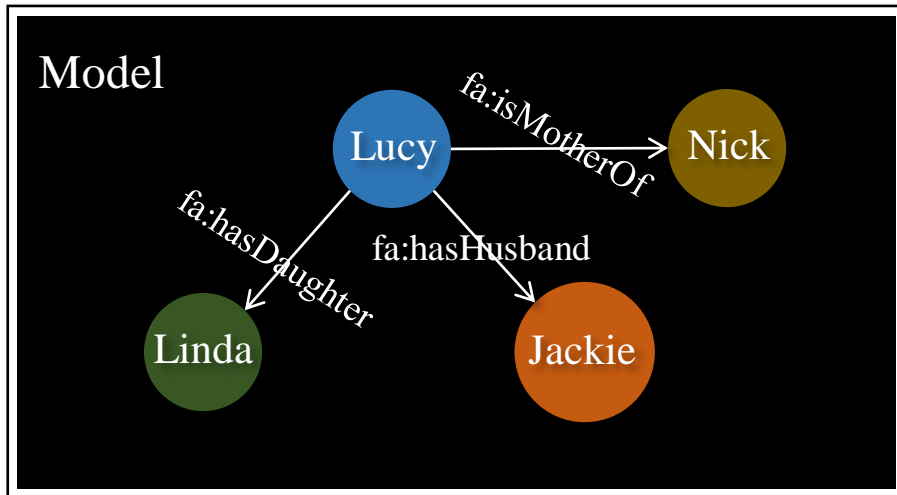


# Jena Inference API

```
[rule1: (?a fa:hasHusband ?b)(?a fa:isMotherOf ?c)->(?b fa:isFatherOf ?c)]  
[rule2: (?a fa:hasHusband ?b)(?a fa:hasDaughter ?c)->(?b fa:isFartherOf ?c)]
```

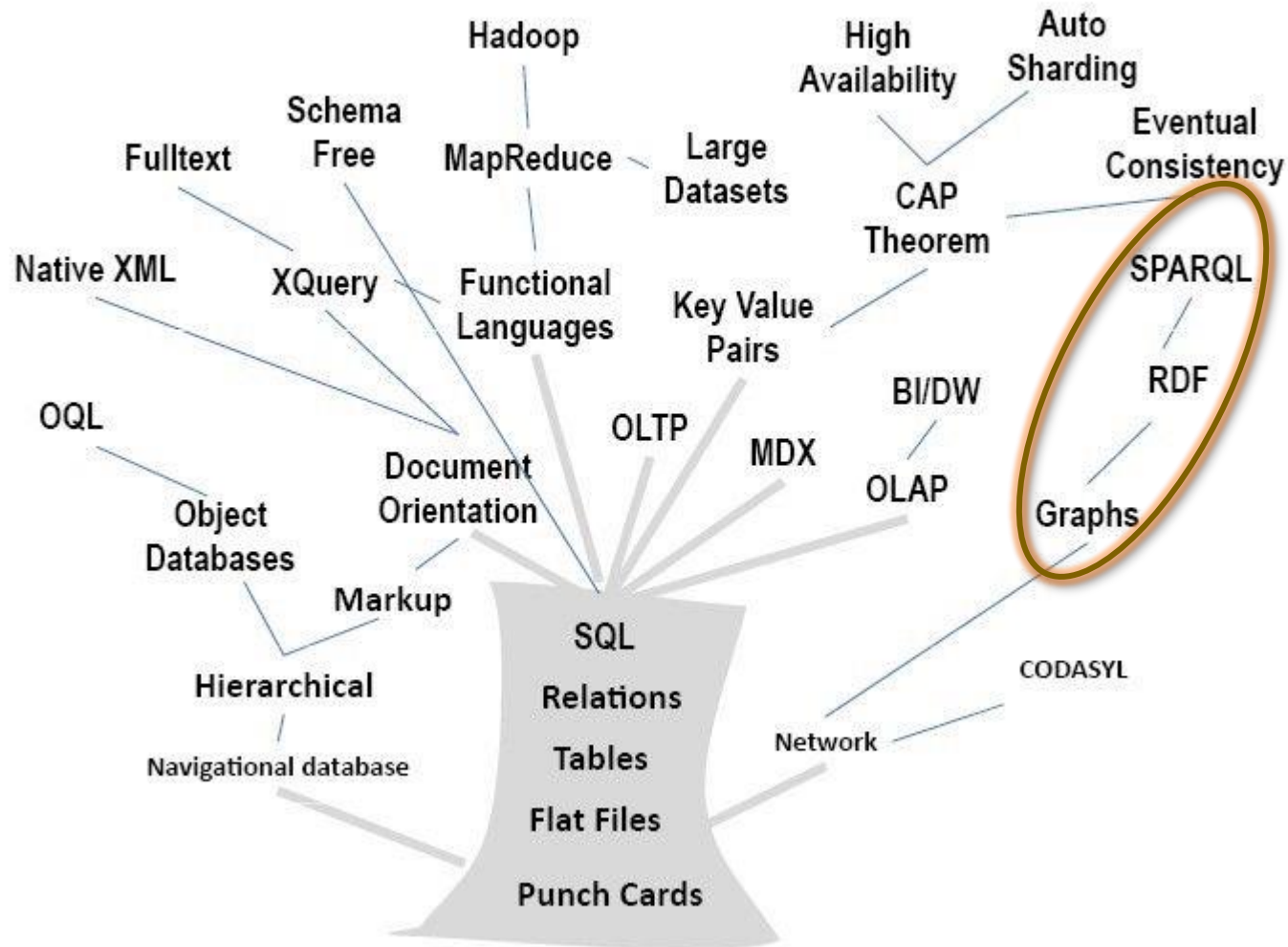
```
Reasoner reasoner = new GenericRuleReasoner(Rule.parseRules(rules));  
InfModel inf = ModelFactory.createInfModel(reasoned, rawData);
```

Lucy fa:hasHusband Jackie ; fa:isMotherOf Nick ; fa:hasDaughter Linda.





# A NoSQL Concept Tree

















## RDF Databases / Triple Store

- **Schema flexibility** - it's possible to do the equivalent of a schema change to an RDF store live, and without any downtime, or redesign .
- **More modern** - RDF stores are typically queried over HTTP, it's very easy to fit them into Service Architectures. Also they handle internationalized content better than typical SQL databases.
- **Standardization** - the level of standardization of implementations using RDF and SPARQL is much higher than SQL. Moving data between stores is easy, as they all speak the same language.
- **Expressivity** - it's much easier to model complex data in RDF than in SQL, and the query language makes it easier to do things like LEFT JOINS (called OPTIONAL in SPARQL).
- **Provenance** - SPARQL lets you track where each piece of information came from, and you can store metadata about it, letting you easily do sophisticated queries.

# RDF Databases / Triple Store

Databases	Customers		
<b>OpenLink Virtuoso</b>			
<b>AllegroGraph</b>			
<b>Stardog</b>			
<b>Oracle NoSQL</b>			

**Others:** Jena TDB, BigData, Sesame, 4store, MarcLogic, YARS2



# How to access Linked Data Sets?

Downloads

Web interface (lookup + services)

SPARQL endpoint

Flint Sparql editor

Virtuoso Faceted Browser



# RDF Serialization

```
model.write(output, "format") ;
```

Jena writer name	RIOT RDFFormat
"TURTLE"	TURTLE
"TTL"	TURTLE
"Turtle"	TURTLE
"N-TRIPLES"	NTRIPLES
"N-TRIPLE"	NTRIPLES
"NT"	NTRIPLES
"RDF/XML-ABBREV"	RDFXML
"RDF/XML"	RDFXML_PLAIN
"N3"	N3
"JSON-LD"	JSONLD
"RDF/JSON"	RDFJSON





# Get the phenotypes of knock-out mouse models for the targets of a selected drug

data sources: DrugBank, HGNC, MGI, BioPortal

Drug: **Imatinib**

## Results:

```
PREFIX dct: <http://purl.org/dc/terms/>
SELECT DISTINCT ?phenotype_label
WHERE {
  SERVICE <http://drugbank.bio2rdf.org/sparql> {
    ?drug <http://bio2rdf.org/drugbank_vocabulary:target> ?target .
    FILTER(?drug = <http://bio2rdf.org/drugbank:DB00619>)
    ?target <http://bio2rdf.org/drugbank_vocabulary:x-hgnc> ?hgnc .
  }
  SERVICE <http://hgnc.bio2rdf.org/sparql> {
    ?hgnc <http://bio2rdf.org/hgnc_vocabulary:x-mgi> ?marker .
  }
  SERVICE <http://mgi.bio2rdf.org/sparql> {
    ?model <http://bio2rdf.org/mgi_vocabulary:marker> ?marker .
    ?model <http://bio2rdf.org/mgi_vocabulary:allele> ?all .
    ?all <http://bio2rdf.org/mgi_vocabulary:allele-attribute> ?allele_type .
    ?model <http://bio2rdf.org/mgi_vocabulary:phenotype> ?phenotypes .
    FILTER (str(?allele_type) = "Null/knockout")
  }
  SERVICE <http://bioportal.bio2rdf.org/sparql> {
    ?phenotypes rdfs:label ?phenotype_label .
  }
}
```

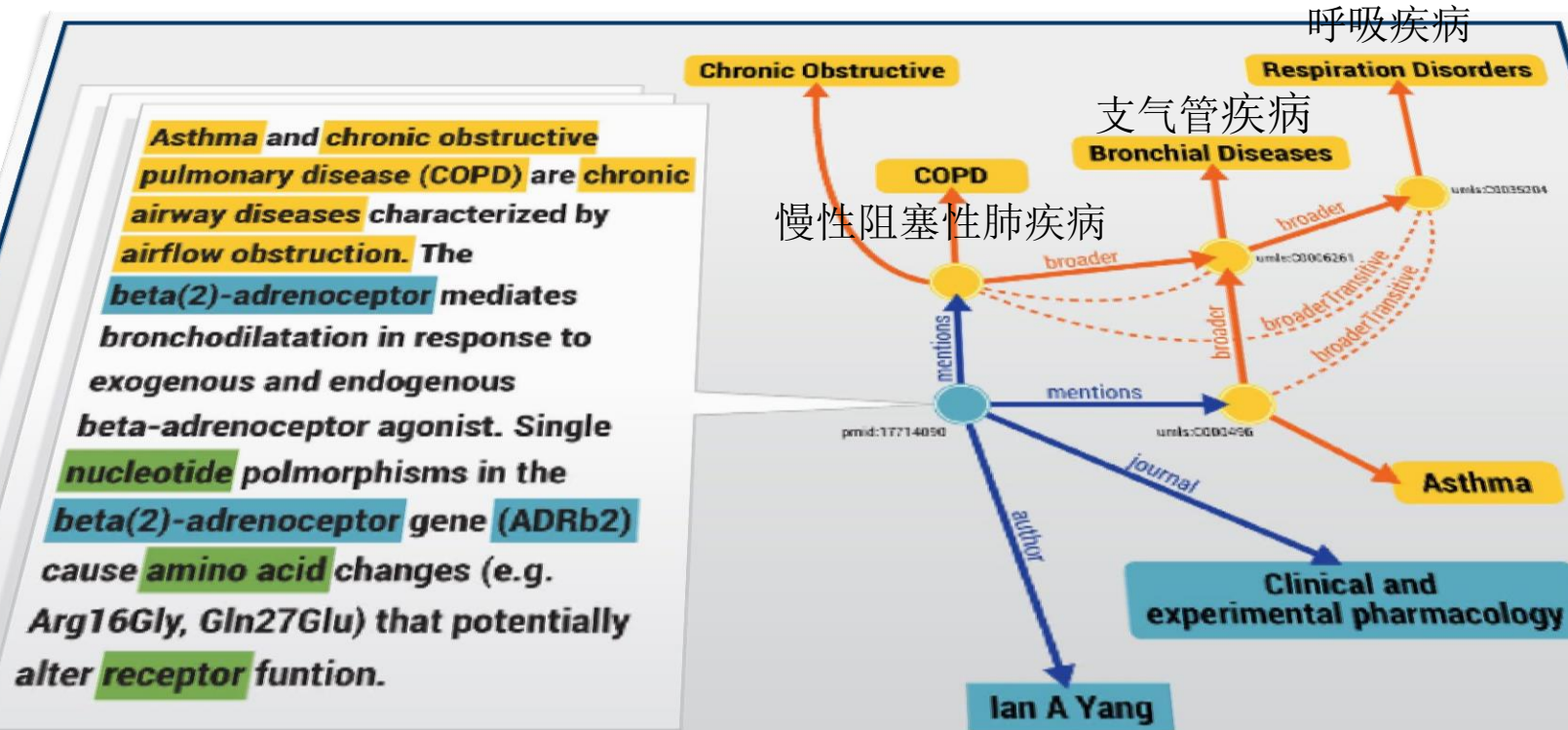
"hemorrhage [mp:0001914]"@en  
"intracranial hemorrhage [mp:0001915]"@en  
"perinatal lethality [mp:0002081]"@en  
"hydrops fetalis [mp:0002192]"@en  
"kidney hemorrhage [mp:0000533]"@en  
"skin hemorrhage [mp:0011514]"@en  
"anemia [mp:0001577]"@en

"出血"  
"颅内出血"  
"围产期的杀伤力"  
"胎儿水肿"  
"肾出血"  
"皮肤出血"  
"贫血"





# Semantic Search

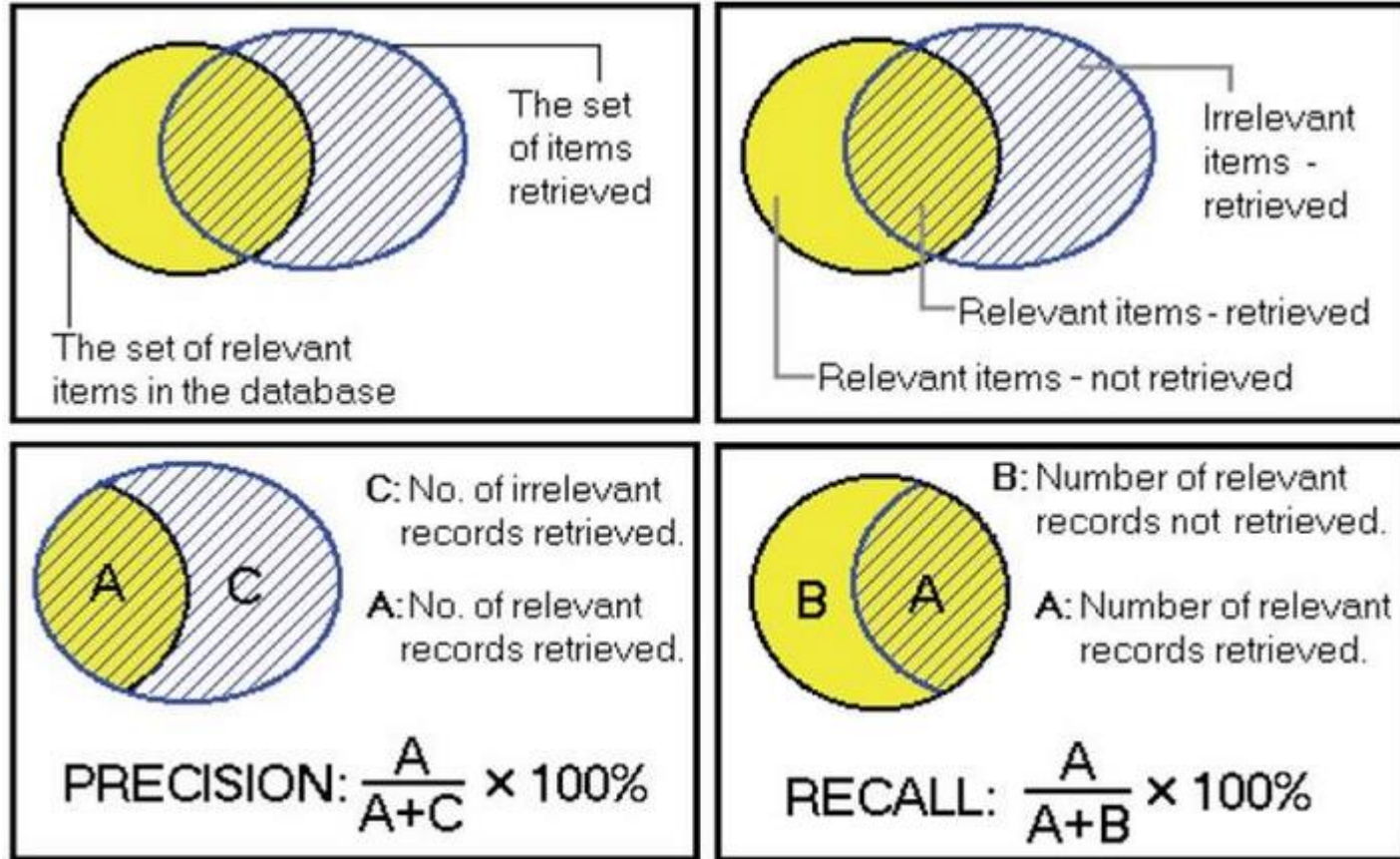


TripleStore

- enriched content
- increased wealth of data
- improved precision and recall



# Precision & Recall

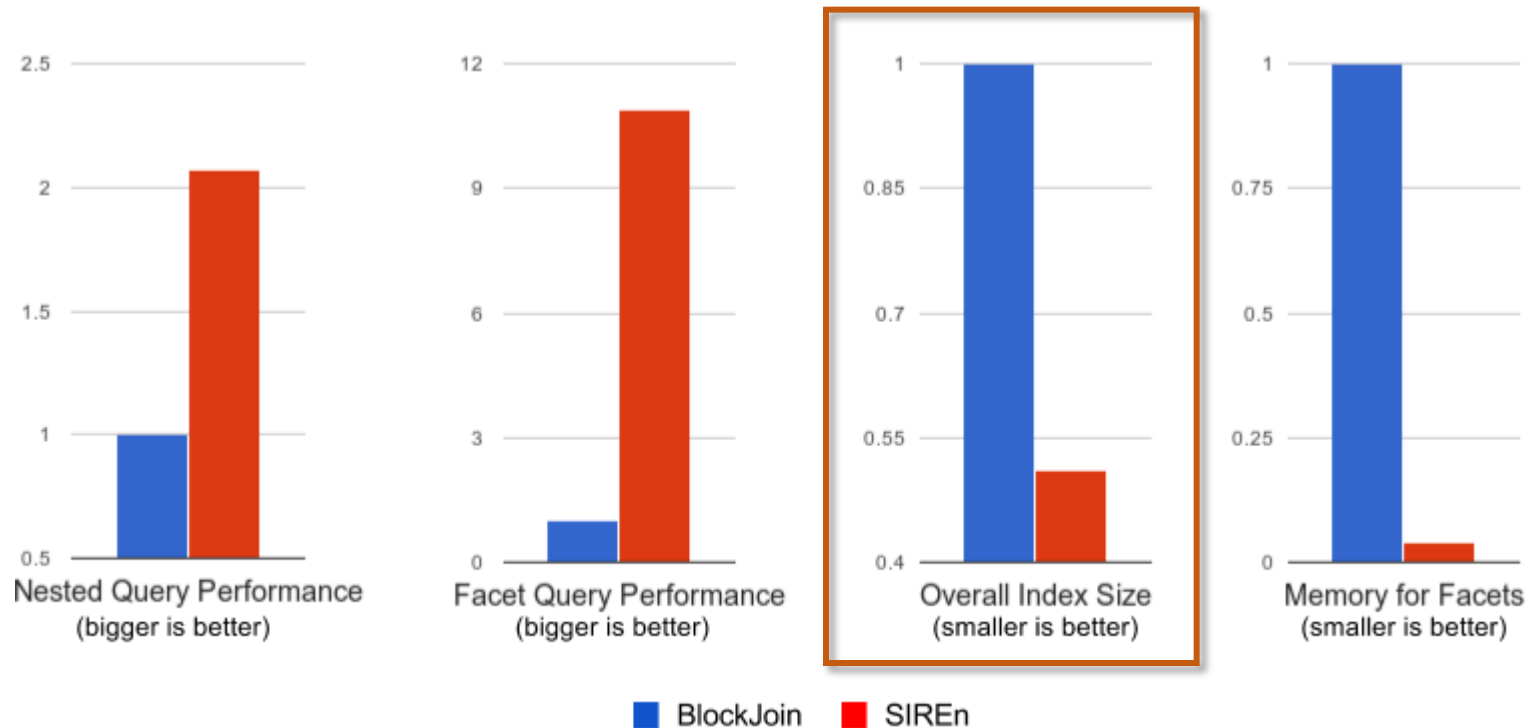




# SIREn {Semantic Information Retrieval Engine}

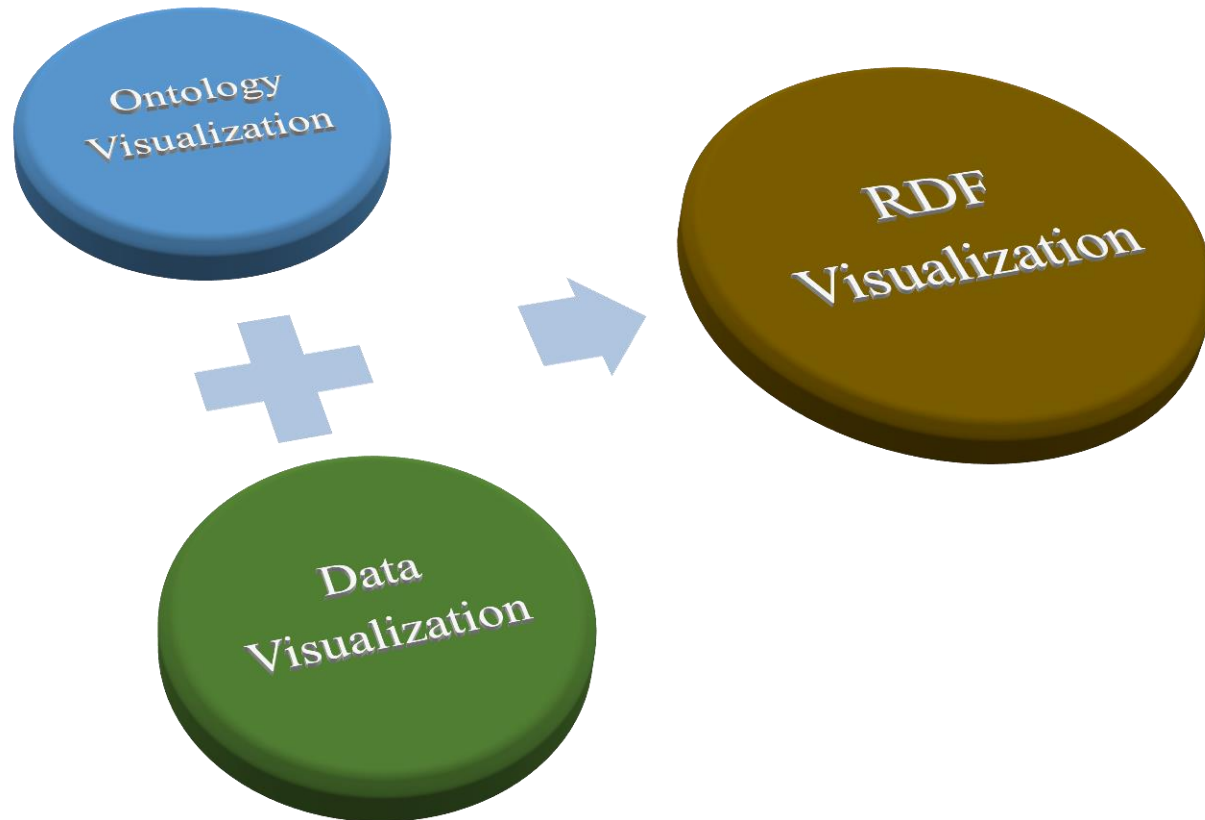
SIREn is a highly scalable open-source full-text search engine especially suited for **nested** and **schemaless** data.

Lucence's BlockJoin:  $100 * 100 = 10,000$  docs  
SIREn:  $100 = 100$  docs





# Visual Tools – RDF Graph





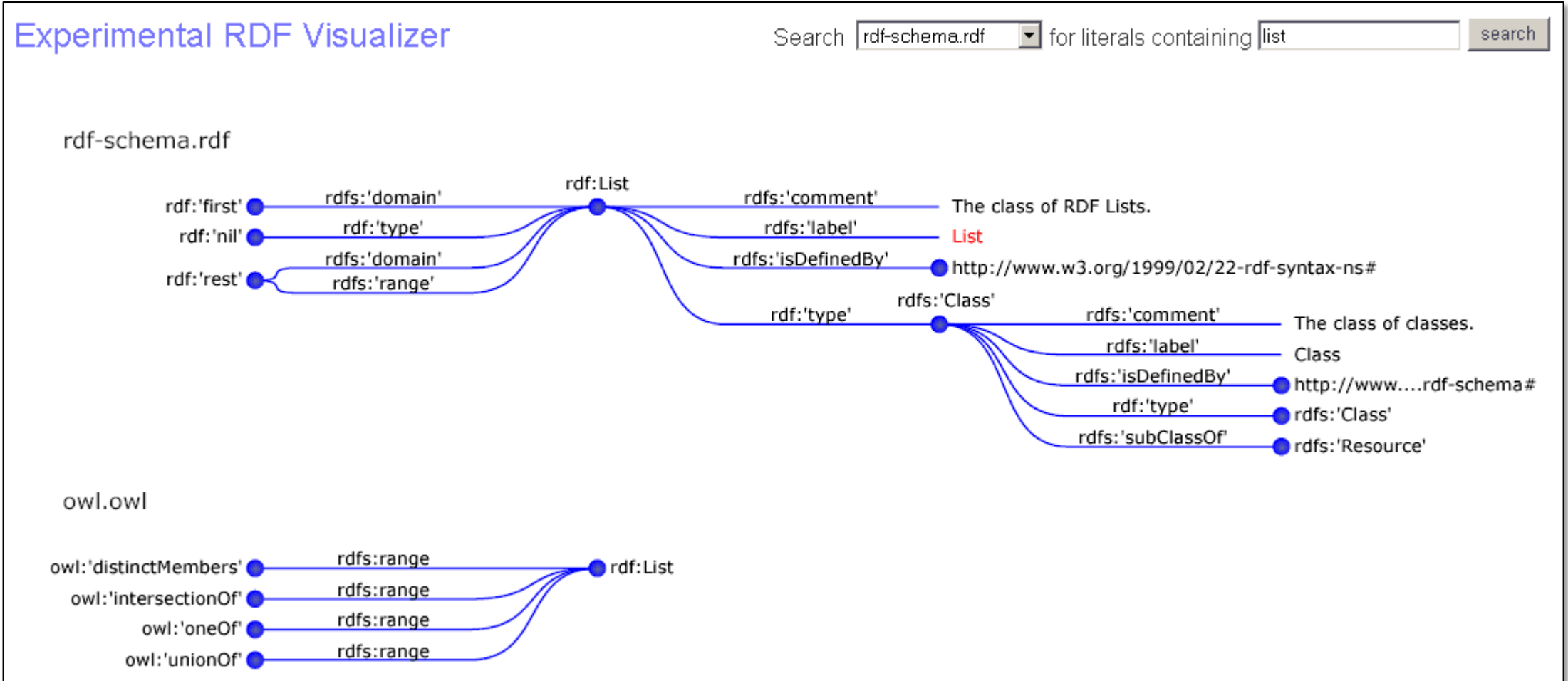
# RDF/OWL Ontology Visualization

- **HP RDF Graph Visualization** describes an experimental node-centric approach to RDF
- **RDF Gravity** is a tool for visualizing RDF/OWL Graphs/ontologies
- **Welkin** is a graph-based RDF visualizer
- **RelFinder** extracts and visualizes relationships between given objects in RDF data and makes these relationships interactively explorable





# HP RDF Graph Visualization





# Welkin

The screenshot shows the Welkin application window. The title bar reads "Welkin". At the top right, there is an "About" button. Below the title bar are "Load" and "Clear" buttons. The main interface is divided into several sections:

- Predicates:** A list of loaded predicates with checkboxes and sliders. The list includes:
  - Predicates
  - <http://www.w3.org> [1127]
    - [/1999/02/22-rdf-syntax-ns](http://www.w3.org/1999/02/22-rdf-syntax-ns) [776]
    - [/2004/02/skos/core](http://www.w3.org/2004/02/skos/core) [351]
  - <http://purl.org> [712]
  - <http://www.imsproject.org> [235]
    - [/rdf/imsmd\\_lifecyclev1p2](http://www.imsproject.org/rdf/imsmd_lifecyclev1p2) [118]
    - [/rdf/imsmd\\_generalv1p2](http://www.imsproject.org/rdf/imsmd_generalv1p2) [117]
  - <http://simile.mit.edu> [106]
    - [/2004/01/ontologies/ocw](http://simile.mit.edu/2004/01/ontologies/ocw) [106]
      - [#linkToImage](#) [106]
- Resources:** A list of loaded resources with checkboxes and sliders. The list includes:
  - Resources
  - <http://ocw.mit.edu>
  - <http://simile.mit.edu>
  - <http://www.imsproject.org> [6]
  - <http://www.w3.org> [1]

The central area displays a network graph with nodes and edges. Several nodes are highlighted with yellow callouts:

- Shigeru Miyagawa**: <http://www.w3.org/2001/vcard-rdf/3.0#FN> -> Shigeru Miyagawa  
<http://www.w3.org/2000/01/rdf-schema#label> -> Shigeru Miyagawa
- John Dower**
- MIT OpenCourseWare**
- Peters, W. T.**
- Visualizing Cultures**

At the bottom left of the graph area, statistics are shown: nodes: 246, edges: 1245, drawing: 52 ms, calculation: 124 ms.

Below the graph are three small plots: "In Degree", "Out Degree", and "Clustering Coefficient".

At the bottom of the window are buttons for "Drawing", "Highlight", and "Parameters". Below these are a "Start" button and a row of checkboxes:  Nodes,  Edges,  Arrows,  Antialias,  Background.



# RelFinder

RelFinder URL 🔍 🔗 🔧

between examples

(1)  ×

(2)  ×

(3)  ×


Filter by: relations: (31/129)

length  class  link  connec...

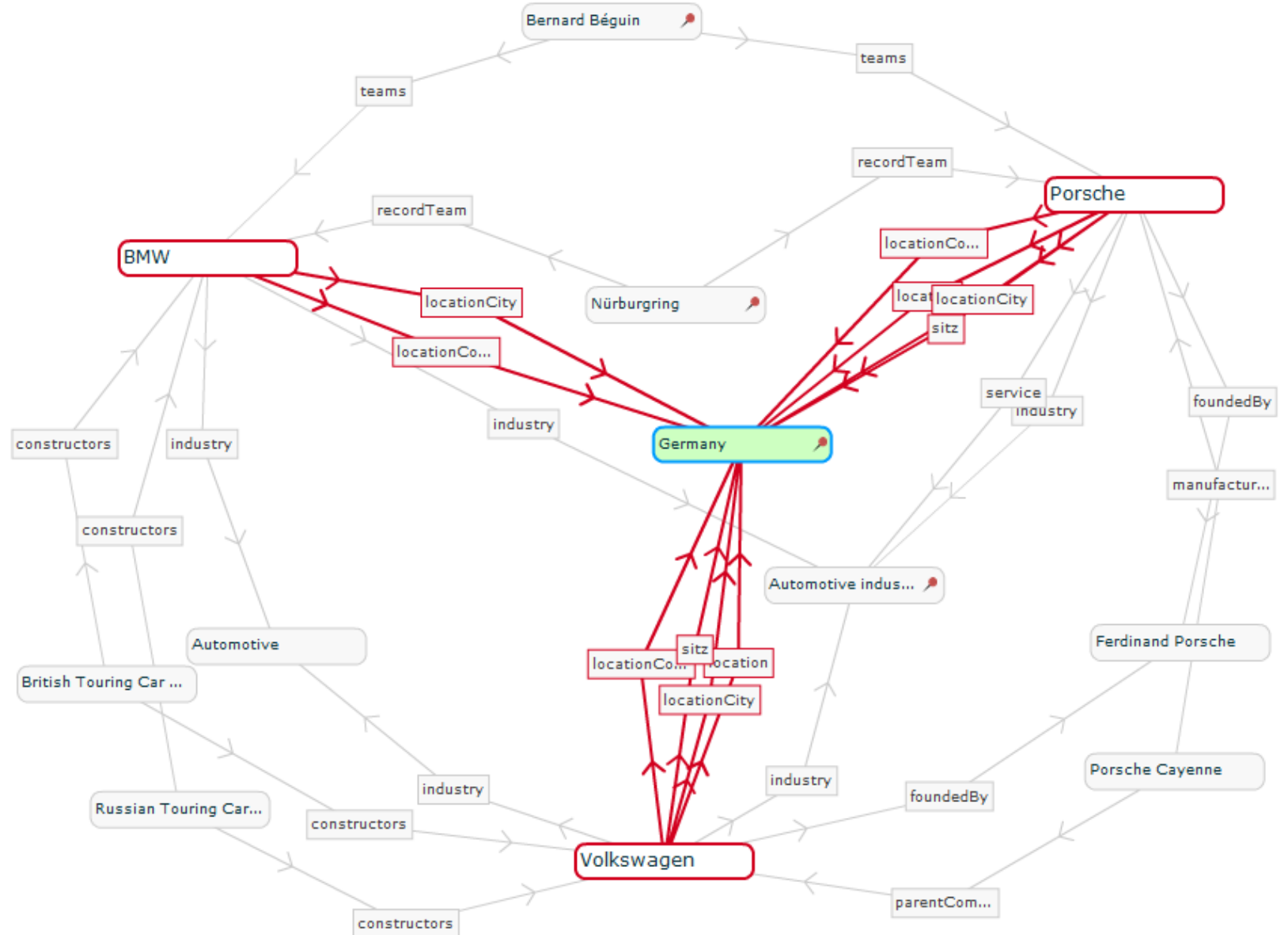
number of objects	num	vi
1	31/34	👁
2	0/95	✖

**Germany** en ▼

More Infos: [dbpedia.org](http://dbpedia.org)  
[www.bundesregierung.de](http://www.bundesregierung.de)



Germany, officially the Federal Republic of Germany, is a federal parliamentary republic in







# RDF Data Visualization

- **RDF-spark** is a graph-based RDF visualizer
- **Sgvizler** provides a javascript based SPARQL query form that passes the SPARQL result set to a powerful visualizer
- **Visual Browser** is a Java application that can visualize the data in RDF scheme
- **Map4rdf** provides your users with a nice map-based visualization of your data





# Sgvizler

- Sgvizler is a javascript which renders the result of SPARQL SELECT queries into charts or html elements.

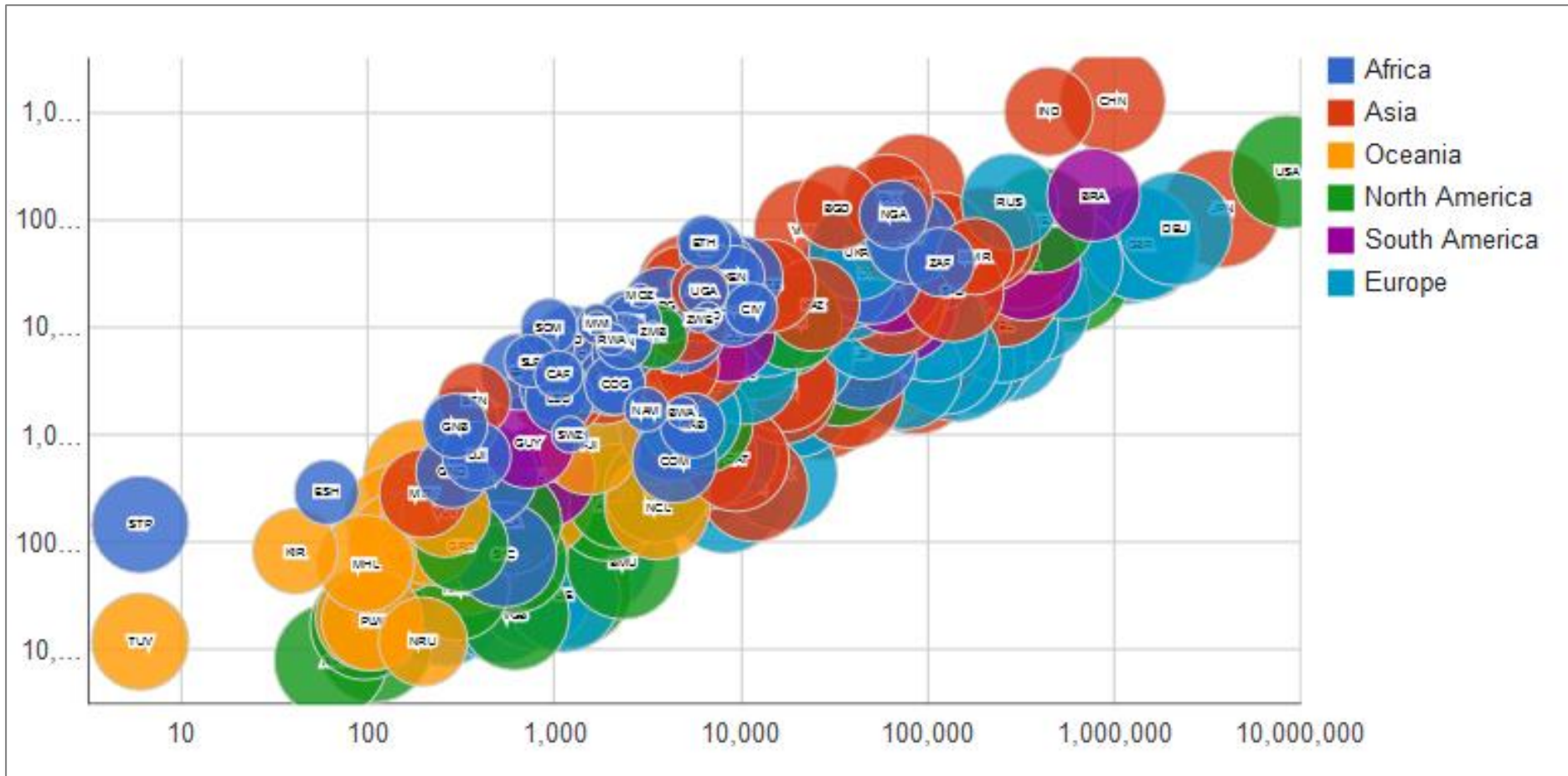
```
<div id="sgvzl_example_query"
  data-sgvizler-endpoint="http://sws.ifi.uio.no/sparql/world"
  data-sgvizler-query="SELECT * WHERE{[] w:hasCountryCode ?ID ; w:hasGNP ?GNP ; w:hasCountryPopulation ?Population ;
w:isCountryInContinent [ rdfs:label ?Continent ] ; w:hasLifeExpectancy ?LifeExpectancy ; }"
  data-sgvizler-chart="google.visualization.BubbleChart"
  data-sgvizler-chart-options="vAxis.logScale=true|hAxis.logScale=true|bubble.textStyle.fontSize=6"
  data-sgvizler-loglevel="2"
  style="width:800px; height:400px;"/>
```

```
SELECT * WHERE {
  [] w:hasCountryCode ?ID ;
  w:hasGNP ?GNP;
  w:hasCountryPopulation ?Population;
  w:isCountryInContinent [ rdfs:label ?Continent ];
  w:hasLifeExpectancy ?LifeExpectancy .
}
```



# Sgvizler

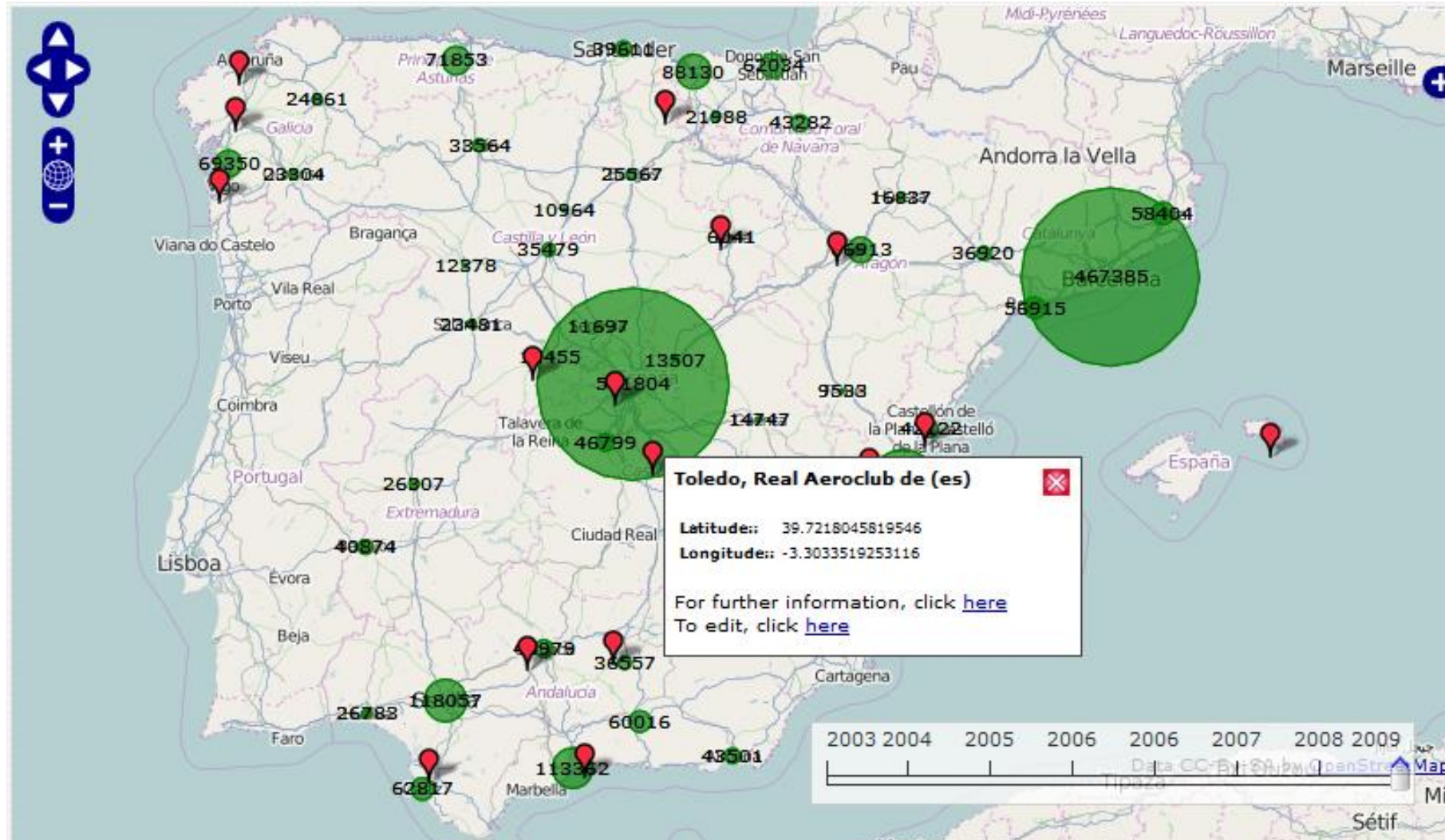
## Result





# Map4Rdf

geo.linkeddata.es





# Map4Rdf

## Toledo, Real Aeroclub de at geo.linkeddata.es

<http://geo.linkeddata.es/resource/Aeroclub/Toledo%2C%20Real%20Aeroclub%20de>



Property	Value
<a href="#">geo:geometry</a>	▪ <a href="http://geo.linkeddata.es/resource/wgs84/39.7218045819546_-3.3033519253116">&lt;http://geo.linkeddata.es/resource/wgs84/39.7218045819546_-3.3033519253116&gt;</a>
<a href="#">rdfs:label</a>	▪ Toledo, Real Aeroclub de (es)
<a href="#">rdf:type</a>	▪ <a href="#">geoes:Aeroclub</a>

### Metadata

Anon\_0

<a href="#">rdf:type</a>	<a href="http://www.w3.org/2004/03/trix/rdfg-1/Graph">&lt;http://www.w3.org/2004/03/trix/rdfg-1/Graph&gt;</a>
<a href="#">foaf:primaryTopic</a>	<a href="http://geo.linkeddata.es/resource/Aeroclub/Toledo%2C%20Real%20Aeroclub%20de">&lt;http://geo.linkeddata.es/resource/Aeroclub/Toledo%2C%20Real%20Aeroclub%20de&gt;</a>
<a href="#">dcterms:creator</a>	<a href="http://geo.linkeddata.es/resource/Organizaci%C3%B3n/InstitutoGeogr%C3%A1ficoNacionalDeEspa%C3%B1a">&lt;http://geo.linkeddata.es/resource/Organizaci%C3%B3n/InstitutoGeogr%C3%A1ficoNacionalDeEspa%C3%B1a&gt;</a>
<a href="#">dcterms:publisher</a>	<a href="http://geo.linkeddata.es">&lt;http://geo.linkeddata.es&gt;</a>
<a href="#">dc:rights</a>	La información geográfica digital comprendida en el Equipamiento Geográfico de Referencia Nacional (artículo 1.1 de la Orden FOM/956/2008) así como los Metadatos de los datos geográficos y servicios del IGN-CNIG, no requieren la aceptación de licencia y su uso será, en cualquier caso, libre y gratuito, siempre que se mencione al Instituto Geográfico Nacional como propietario de los datos.
<a href="#">dcterms:spatial</a>	<a href="http://geo.linkeddata.es/resource/Pa%C3%ADs/Espa%C3%B1a">&lt;http://geo.linkeddata.es/resource/Pa%C3%ADs/Espa%C3%B1a&gt;</a>
<a href="#">prv:createdBy</a>	Anon_1 ( <a href="#">more</a> )

[expand all](#)

This page shows information obtained from the SPARQL endpoint at <http://geo.linkeddata.es/sparql>.

[As Turtle](#) | [As RDF/XML](#) | [Browse in Disco](#) | [Browse in Tabulator](#) | [Browse in OpenLink Browser](#)



# Visual Tools

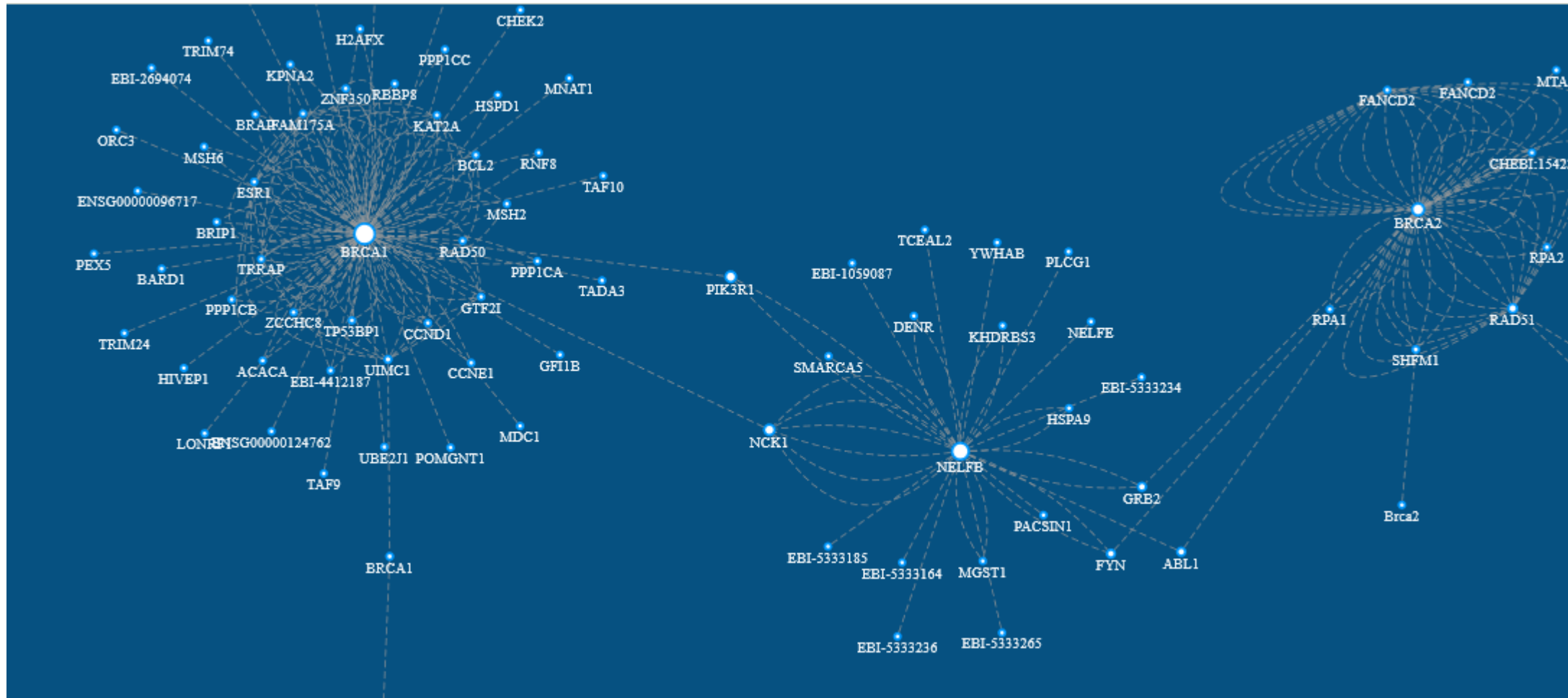
- **D3**
- **Simile Exhibit**
- **Echarts**
- **Google Visualization API**
- **Data Publica**
- **GeoAPI**
- **Cytoscape**





# Cytoscape

*An open source software platform for visualizing complex networks and integrating these with any type of attribute data.*





# BIO-Linked Scientific Data

会员分布

141 ↑ Member



查询

研究领域

- 76 分子生物学
- 8 基础医学
- 5 免疫学
- 1 人类学
- 5 神经生物学
- 4 生理学
- 5 生物工程
- 78 生物化学
- 3 生物学其他学科
- 2 微生物学
- 4 细胞生物学
- 7 药理学
- 8 遗传学
- 2 肿瘤生物学

地域

- 1 澳门
- 1 包头市
- 39 北京市
- 3 成都市
- 2 重庆市
- 1 大连市





# VISUalization Playground

## Linked Open Aalto Data Service

---

### VISUalization Playground

Add endpoints + -

<http://data.aalto.fi/sparql>

Class usage in certain domain

Show available prefixes

```
SELECT ?class (count(?s) as ?count) WHERE { ?s a ?class } GROUP BY ?class
```

Prefixify uris

Visualize



Aalto University

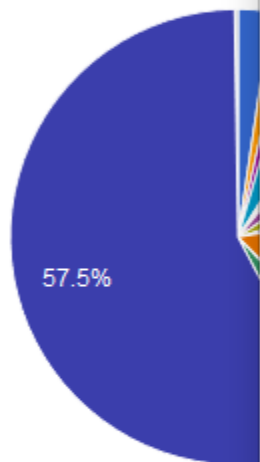


# Linked Open Aalto Data Service

## VISUalization Playground

Visualize Edit Export

VISUalisation



### Chart Editor

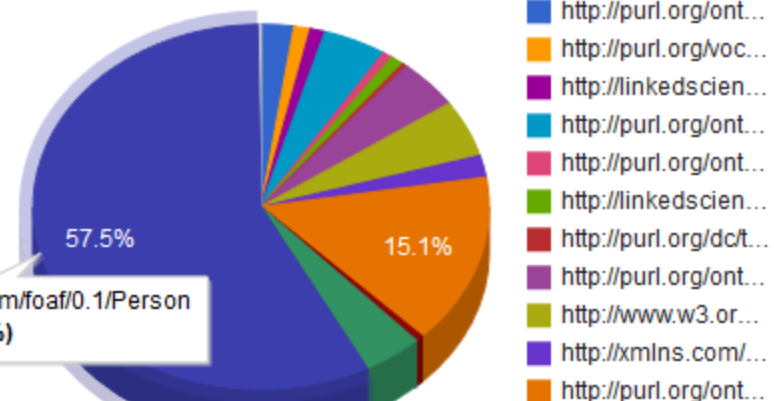
Start Charts Customize

Chart name

- Line
- Area
- Column
- Bar
- Scatter
- Pie**
- Map
- Trend
- More

VISUalisation

<http://xmlns.com/foaf/0.1/Person>  
225725 (57.5%)



▲ 1/2 ▼

OK

Cancel



# R script

R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, ...) and graphical techniques.

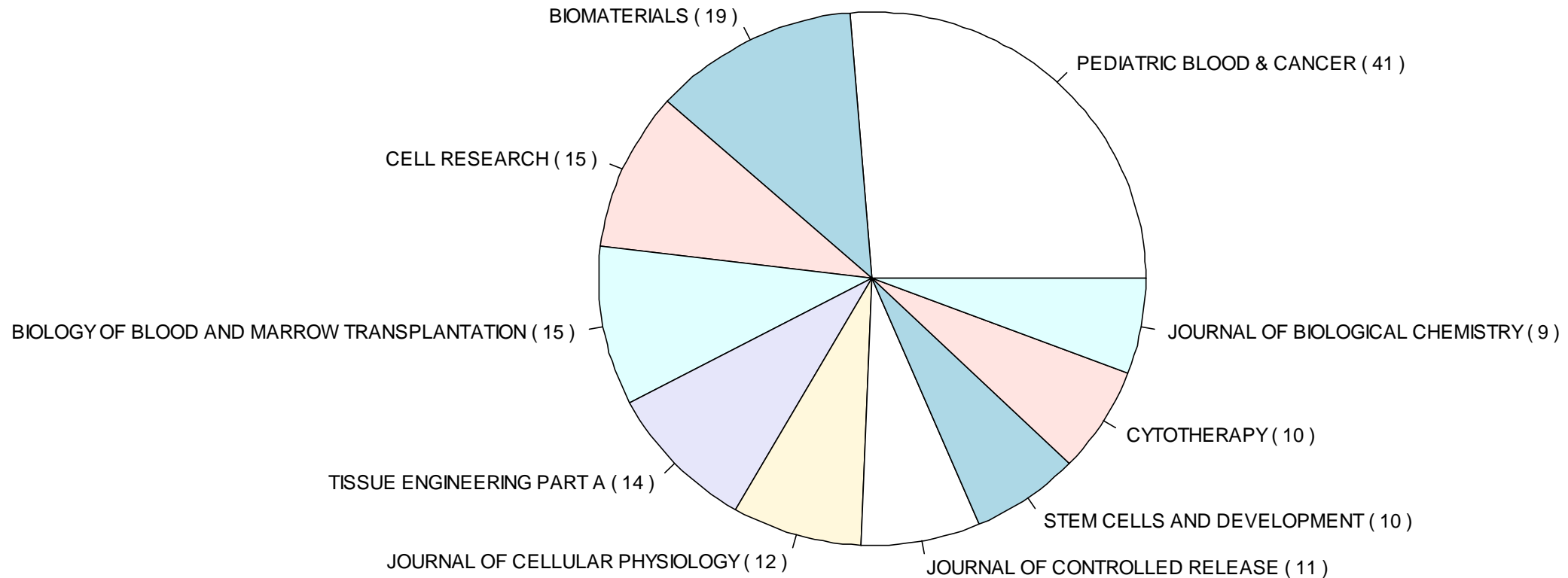


```
1 library(SPARQL)
2 library(igraph)
3 library(network)
4 library(ergm)
5
6 # Live DBpedia endpoint
7 endpoint <- 'http://localhost:8890/sparql'
8 options <- NULL
9
10 prefix <- c("db", "http://dbpedia.org/resource/")
11
12 sparql_prefix <- "PREFIX bibo: <http://purl.org/ontology/bibo/>
13                   PREFIX dc: <http://purl.org/dc/elements/1.1/>
14                   PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
15                   PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
16 "
17
18 q <- paste(sparql_prefix,
19           'select ?journal (count(distinct ?s) as ?count)
20           where {graph <http://www.cba.ac.cn/stemcell/graph/article> {
21             ?s dc:title ?title ; bibo:presentedAt ?journal .}
22           } order by desc(?count) limit 10')
23
24 res <- SPARQL(url=endpoint,query=q,ns=prefix)$results
25 # iconv(res,"utf-8","gbk")
26
27 slices <- res$count
28 lbls <- res$journal
29 lals <- paste(lbls, slices)
30 pie(slices,labels=lbls,main="Pie Chart of Journal")
```



# R - Pie Chart

**Pie Chart of Journal**



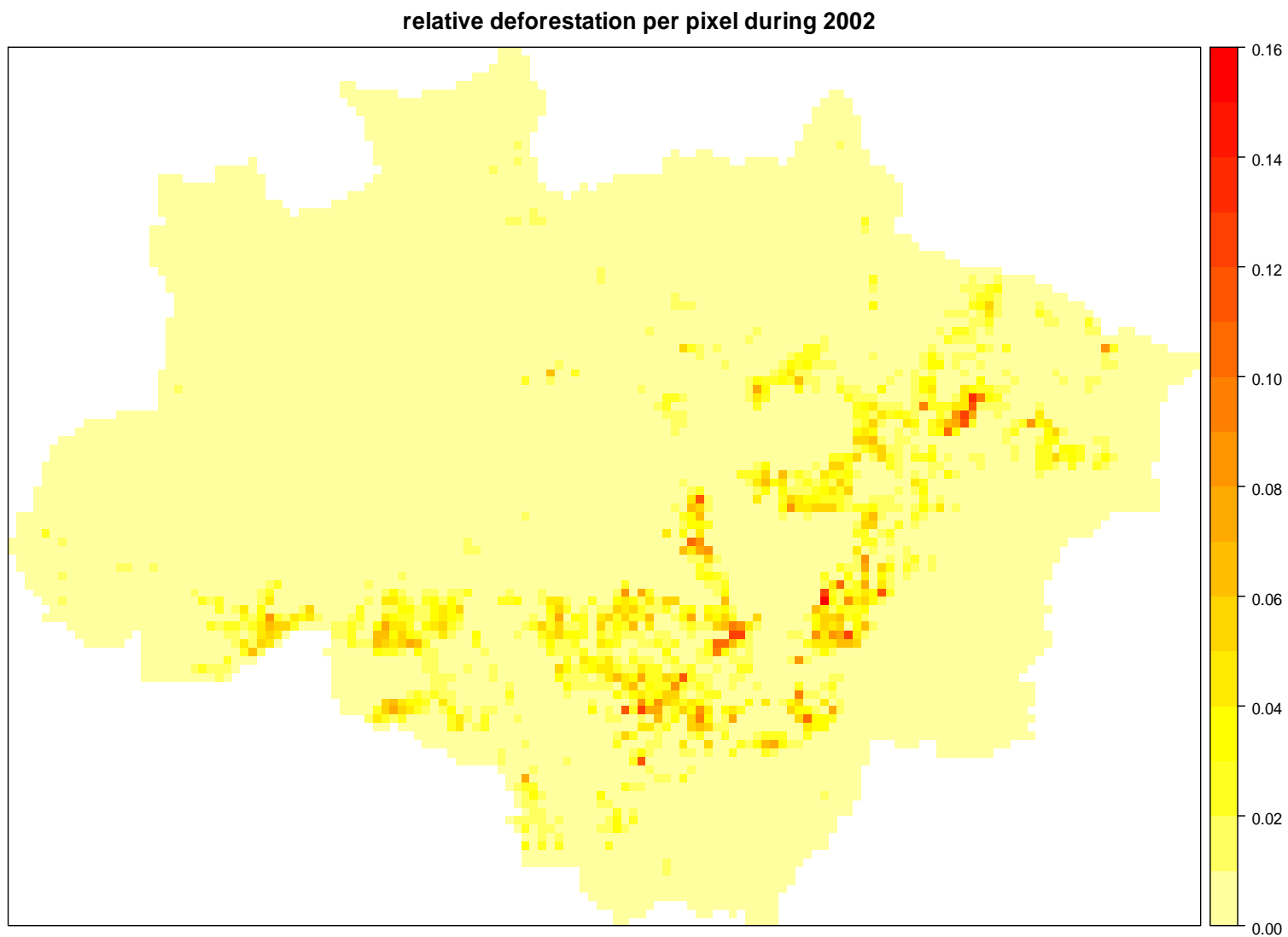


# R - Brazilian Amazon Rainforest

```
1 library(SPARQL)
2 library(sp)
3 endpoint <- "http://spatial.linkedscience.org/sparql"
4 q <- "SELECT ?cell ?row ?col ?polygon
5 WHERE {
6   ?cell a <http://linkedscience.org/lsv/ns#Item> ;
7   <http://spatial.linkedscience.org/context/amazon/Lin> ?row ;
8   <http://spatial.linkedscience.org/context/amazon/Col> ?col ;
9   <http://observedchange.com/tisc/ns#geometry> ?polygon .
10  }"
11
12 res <- SPARQL(url=endpoint, q)$results
13 for(var in c("DEFOR_2002", "DEFOR_2003", "DEFOR_2004", "DEFOR_2005", "DEFOR_2006",
14   "DEFOR_2007", "DEFOR_2008")) {
15   tmp_q <- paste("SELECT ?cell ?", var, "\n
16   WHERE { \n ?cell a <http://linkedscience.org/lsv/ns#Item> ;\n
17   <http://spatial.linkedscience.org/context/amazon/", var, "> ?", var, " .\n } \n", sep="")
18   cat(tmp_q)
19   res <- merge(res, SPARQL(endpoint, tmp_q)$results, by="cell")
20 }
21 amazon <- res
22 amazon$row <- -res$row
23 coordinates(amazon) <- ~ col+row
24 gridded(amazon) <- TRUE
25 spplot(amazon, "DEFOR_2002", col.regions=rev(heat.colors(17))[-1], at=(0:16)/100,
26   main="relative deforestation per pixel during 2002")
27 spplot(amazon, c("DEFOR_2002", "DEFOR_2003", "DEFOR_2004", "DEFOR_2005",
28   "DEFOR_2006", "DEFOR_2007", "DEFOR_2008"),
29   col.regions=rev(heat.colors(26))[-1], at=(0:20)/80, as.table=T,
30   main="relative deforestation per pixel")
```

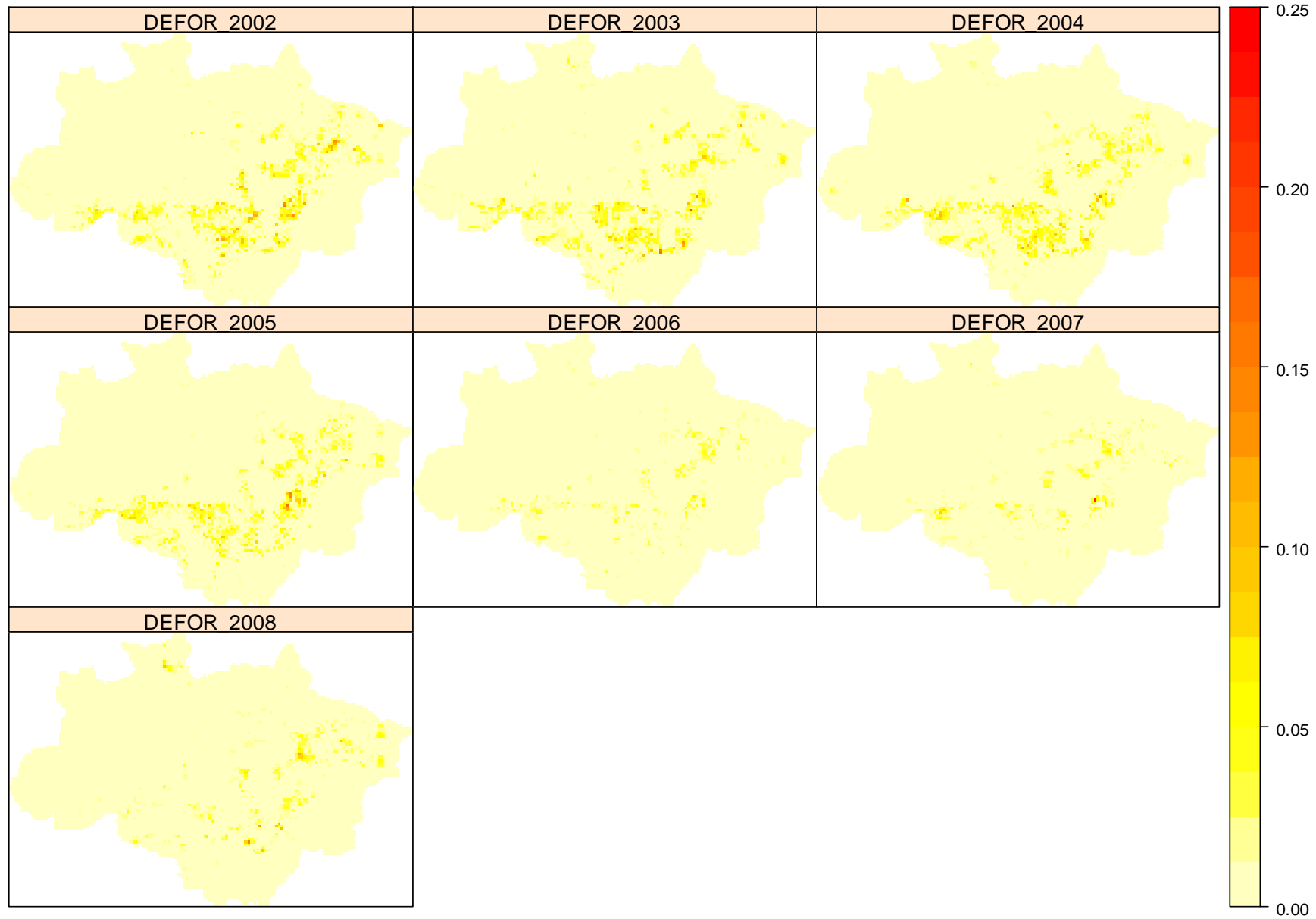


# R - Brazilian Amazon Rainforest (2002)





relative deforestation per pixel





# OpenLink Virtuoso Open-Source Edition

- **Install**
- **Use it for linked data**
- **Install VAD plugin packages**
- **Visualization interface (Demo datasets: Map, Article)**
  - **Sparql endpoint**
  - **iSparql**
  - **Facet view**
  - **Pivot viewer**





**Give a little, Get a little!**

