# Theoretical part
# Semantic Web

Olivier DAMERON -
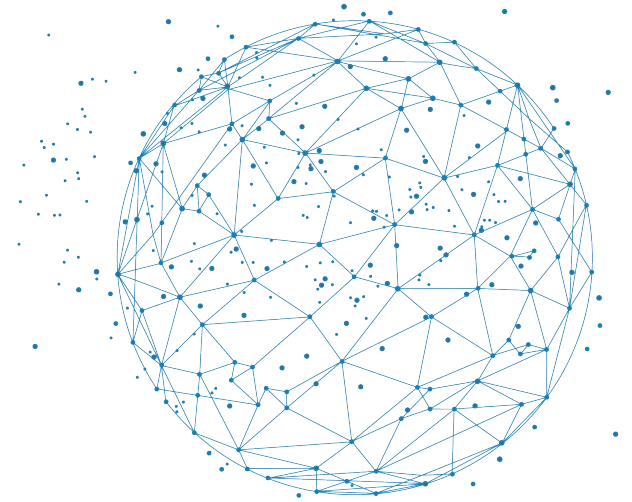Alban GAIGNARD -
Pierre LARMANDE

- Life science data require
  - Integration
  - Knowledge-based reasoning
- The Semantic Web provide a relevant framework
- Use RDF to represent knowledge graphs
- Use SPARQL to query knowledge graphs
- Use RDFS and OWL to formalize knowledge as ontologies

What you will learn (hopefully):

- A general understanding of metadata and (symbolic) knowledge…
- … that relies on surprisingly simple principles

# Life science data from an information science perspectives

**Big data** [Laney, 2001]

Datasets so **large** and **complex** that traditional data processing is inadequate

**Life science: data deluge since the 90s** [Aldhous, 1993]

- Computerized biomedical data (evidence-based medicine, translational medicine, precision medicine)
- 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]

Nat Rev Drug Discov. 2002 Jun;1(6):479.

**Dealing with the data deluge.**

Lanfear J[1].

⊕ **Author information**

PMID: 12119750 [PubMed - indexed for MEDLINE]

Our estimation is that genomics is a ``four-headed beast'' – it is either
on par with or the most demanding domain[...] in terms of:

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

PERSPECTIVE

# 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

5

- multiple scale (heterogeneity)
- (highly) interdependent at each scale
- interdependent between scales
- variability
- incompleteness
- evolutive
- distributed
  (and lack of interoperability)

# Challenge (computational): How to handle this complexity?

- Experts are very good at analyzing data on their domain (hint…)
  - on their domain
  - on their data
  - with their limited human capacity
  - … help is (badly) needed!
- The difficulty is to analyze data systematically
- Expertise = ability to use knowledge for interpreting data
  - aggregating data
  - inferring connections
- How to use their expertise, instead of trying to re-discover it?
  - assumption: we may not always need to re-discover naively the whole biology at each experiment

- *Requirement 1:* identify resources with interoperable identifiers
- *Requirement 2:* describe resources
    - their characteristics
      (e.g. start and end position of a gene)
    - their relations to other resources
      (e.g. the transcript associated to a gene, the transcription factors that regulate it…)
    - the categories they belong to
- *Requirement 3:* combine descriptions from different origins
- *Requirement 4:* query these descriptions
- *Requirement 5:* support semantically-rich querying and reasoning (because of the inner complexity) using domain knowledge

``Metadata, you see, is really a love note - it might be to yourself, but in fact it's a love note to the person after you, or the machine after you, where you've saved someone that amount of time to find something by telling them what this thing is."

Jason Scott - http://ascii.textfiles.com/archives/3181

Annotation = explicit representation of the result of some interpretation process

Annotation = explicit representation of the result of some interpretation process

- ideally by an expert (from big data to smart data)
- requires some background knowledge
- formalization spanning the whole semantic spectrum, ranging
  - from free text…
  - …to controlled vocabularies…
  - …to (shared) semantic frameworks

Annotate data =

- Describe **explicitly**…
  - … the **relevant elements** in your data…
  - … their **characteristics**…
  - … and the **relations between them**
- So that users (you + the non-experts) or programs do not have to go once again through the (tedious, complicated) process of interpreting them
  - It is important
  - It seems easy…
    - *Exercice: describe a set of images*
  - …but turns out more difficult than expected
- 2 aspects:
  - Interpreting and describing data is for domain experts
  - Formalizing and representing the annotations and their dependencies is for data engineers

*Data / Observations*

**#1** Annotate data = describe their interpretation -> metadata

**#2** There are some dependencies btw some annotations -> knowledge graph

*Abstractions / Models*

**#3** The dependencies between categories support multiple layers of generalization -> ontologies

*Graph of data U knowledge graph U ontologies*

**#4** Reasoning = rules for traversing the graphs

*Add annotations? But we have too much data already!*

Benefits

- Can be used as proxy to complex data
- Simplifies by providing a compact abstraction
- Overcomes variability
- Enriches by making explicit the underlying meaning

**Storing, sharing and reusing these annotations is the key to life science systematic data analysis**

# General introduction to semantic web

# Knowledge Graphs

*Several facets of Artificial Intelligence*

**Named entity recognition (NLP)**
➜ « pasteur » refers to Louis Pasteur

**Knowledge Representation**
➜ Louis Pasteur is a scientist
➜ a scientist is a kind of person

**Reasoning**
➜ « if an entity is a scientists then its also a person »
➜ « If a web resource is a person then display his/her portrait »

**… and many other such as Computer Vision, Machine Learning**

1

Knowledge representation and reasoning : still needed for question answering systems !

# Siri - Software Engineer - Knowledge Graph

Santa Clara Valley (Cupertino), California, United States
Machine Learning and AI

☆  ⬆

Envoyer un CV

‹ Retour aux résultats de recherche

## Summary

Posted: **May 29, 2019**
Weekly Hours: **40**
Role Number: **200039719**

The Knowledge Graph team is looking for outstanding engineers to build the next- generation of knowledge graph and data infrastructure at Apple to power features including Siri and Spotlight. If you are interested in building a world-class

## Key Qualifications

- Extensive systems programming experience in either Python or Java. Solid system development skills in UNIX-type OS (e.g. Linux, Mac OS)
- Experience working with large data sets and pipelines, ideally using the Apache
- software stack (e.g. Spark, HBase)
- Excellent problem-solving and analytic skills
- You are self-motivated and able to quickly learn new domains
- You have good attention to detail
- Broad knowledge of computer science and systems
- Excellent communication and collaborative skills; Able to work as part of a small, focused team and give your best effort

## Description

The Siri Knowledge Graph team is building groundbreaking technology in the areas of question answering, knowledge base construction and machine learning. We aim to be a "know-it-all" question answering system, capable of answering questions from hundreds of millions of users about nearly anything. The question answering system is backed by a knowledge graph that was automatically constructed from a vast number of data sources including natural language text, HTML tables, and many others. You will have exciting opportunities to working on rapidly building a more complete and accurate knowledge graph with impact across all of Apple. This knowledge graph also enables many other features across Apple besides the question answering feature at Siri.
The problems we pursue include

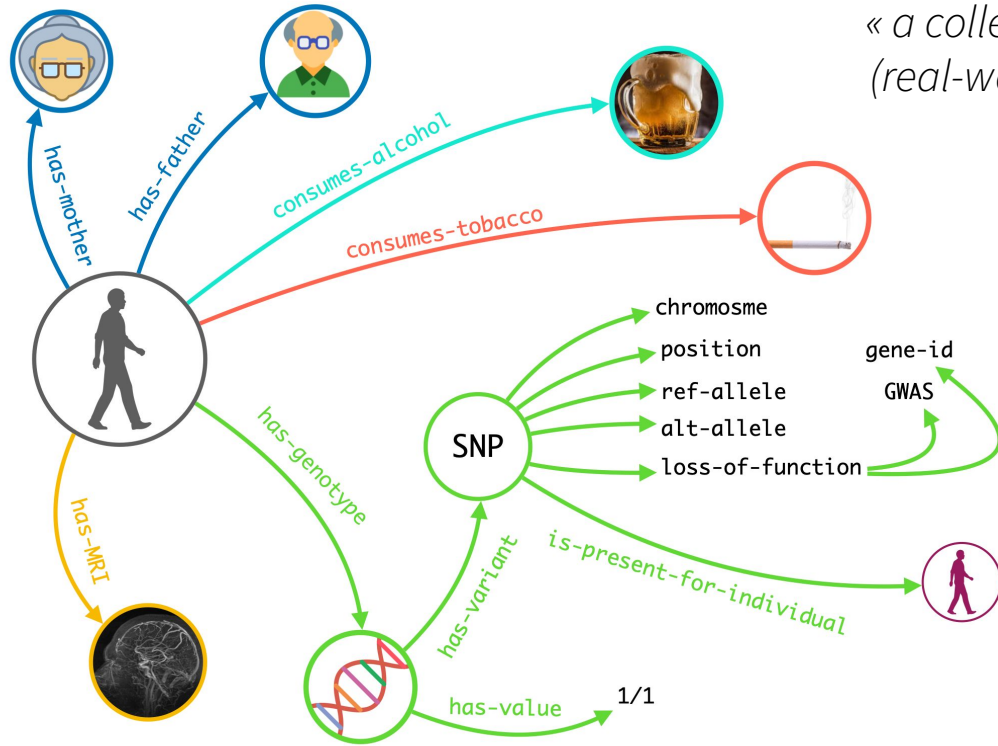* Information Extraction from natural language text and semi-structured data such as HTML tables
* Data Integration (e.g., Entity Resolution and Knowledge Fusion)
* Knowledge graph reasoning and inference
* Named Entity Linking

1

# What is a Knowledge Graph ?



*« a collection of **interlinked descriptions of things** (real-word objects, abstract concepts, events, etc.) »*

a **database**
→ information storage / extraction
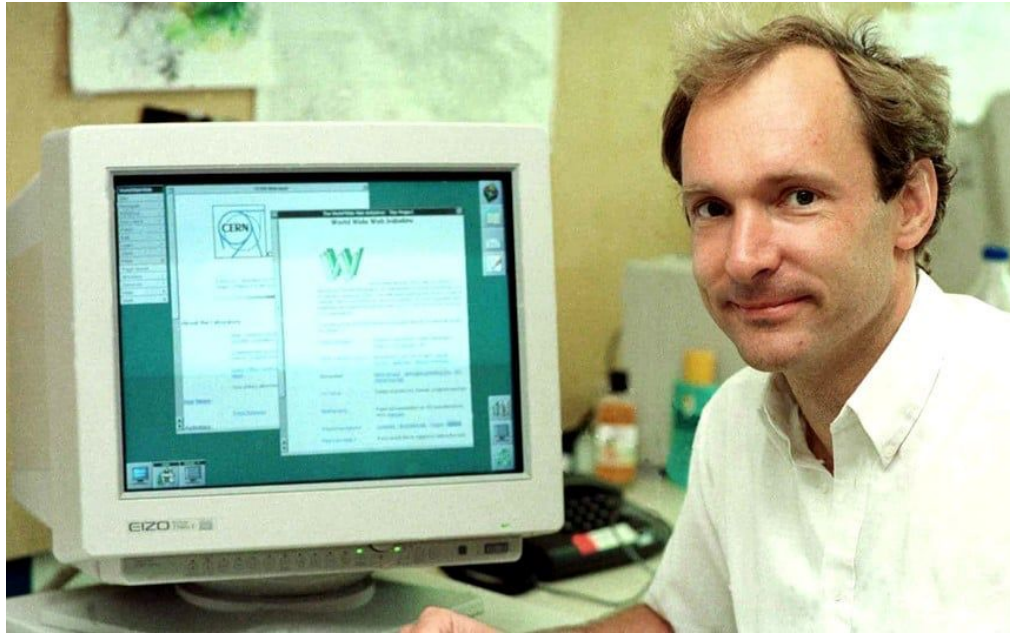
a **graph**
→ network analysis

a **knowledge base**
→ formal semantics
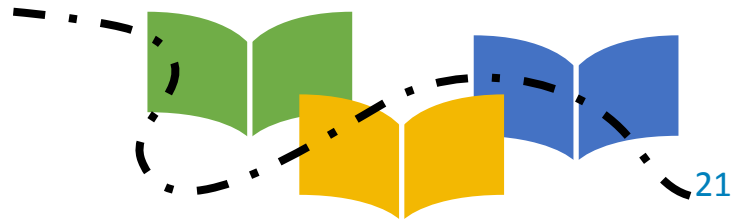(logical facts, logical inferences)

From **linked documents** (Web) …

… to **linked data** (Semantic Web)

A **de-centralized** system of **hypertext documents** based on URL, HTTP, and HTML **standards** (World Wide Web consortium, W3C).

**Surface web < 10% = indexed by search engines**

Bing

Google

Wikipedia

...

Youtube

...

Scientific papers

Academic databases

...

Medical records

Social media

Legal documents

**Deep web > 90%**

# Why so much data is not findable ?





## Genetic sequences

- 1st line = label

- 2nd line = raw sequence (A,T,C,G)

- 4th line = quality score for each base

## Medical image

- MRI as imaging modality

- brain as imaging target

- labeled data (intensity/volume ?)

Make  and  better exchange, interpret & reason on diverse data !!

```
@HWI-ST534_129:2:24:20503:16510:CGATGT
CTGAGAGCCGGGGAAGCCGGCGGAGCCGGGGGACTGGCGAGCCGGAACAT
+
HHHHHHHHHEFDDGDDFBFGG>7D4<9;<&?:;<DC>CCDD@?=?A###
@HWI-ST534_129:2:42:2118:9580:CGATGT
GGCGGAGCCGGGTGACTGGCGAGCCGGAACATCAGGCGCCGCCGCAGAGA
+
GEECGGGBGIDF6FFFFEF=IDEFBEE8E8E?EEB@6=9B##########
@HWI-ST534_129:2:2:12654:80229:CGATGT
CGGAGCCGGGTGACTGGCGAGCCGGAACATCAGGCGCCGCCGCAGAGAAG
+
GGEGFCDCBBAEEEEGGFGFG;EGEEGFFBBEDEBDFGFCFF;DF2D<DD
@HWI-ST534_129:2:48:12356:179714:CGATGT
GAGCCGGGTGACTGGCGAGCCGGAACATCAGGCGCCGCCGCAGAGAAGAA
+
E=GHFHEGHHBCGDDBEEBBCBDDDE@EGBD=ABDCB?EC;@@8@EEB;E
@HWI-ST534_129:2:44:8225:39540:CGATGT
GGGTGACTGGCGAGCCGGAACATCAGGCGCCGCCGCAGAGAAGAACTATG
+
HHHEHHHHHHHHGHHHHFHHHHHHHHDFDHHBHFDDFEFEFF>G<CCCE
```

# W3C recommendations



http://dbpedia.org/resource/RAC1

.csv

abstract

Rac1, also known as Ras-related C3 botulinum toxin substrate 1, is a protein found in human cells. It is encoded by the RAC1 gene. This gene can produce a variety of alternatively spliced versions of the Rac1 protein, which appear to carry out different functions. (en)

RAC1

→ **technological stack to implement FAIR principles** 😀

# Wikipedia … for humans

# Wikipedia ... for machines (DBpedia)

```
SELECT DISTINCT ?gene ?entrez_id ?uniprot_id WHERE {
    ?gene dbo:abstract ?abstract .
    FILTER (regex(?abstract, "toxin")).
    ?gene dbo:entrezgene ?entrez_id .
    OPTIONAL {?gene dbo:uniprot ?uniprot_id} .
}
```

http://dbpedia.org/sparql

Virtuoso SPARQL Query Editor
About | Namespace Prefixes | Inference rules | RDF views | iSPARQL

Default Data Set Name (Graph IRI)
http://dbpedia.org

Query Text

(Security restrictions of this server do not allow you to retrieve remote RDF data, see details.)
Results Format:          HTML
Execution timeout:       30000          milliseconds (values less than 1000 are ignored)
Options:                 ☑ Strict checking of void variables
                         ☐ Log debug info at the end of output (has no effect on some queries and output formats)
                         ☐ Generate SPARQL compilation report (instead of executing the query)
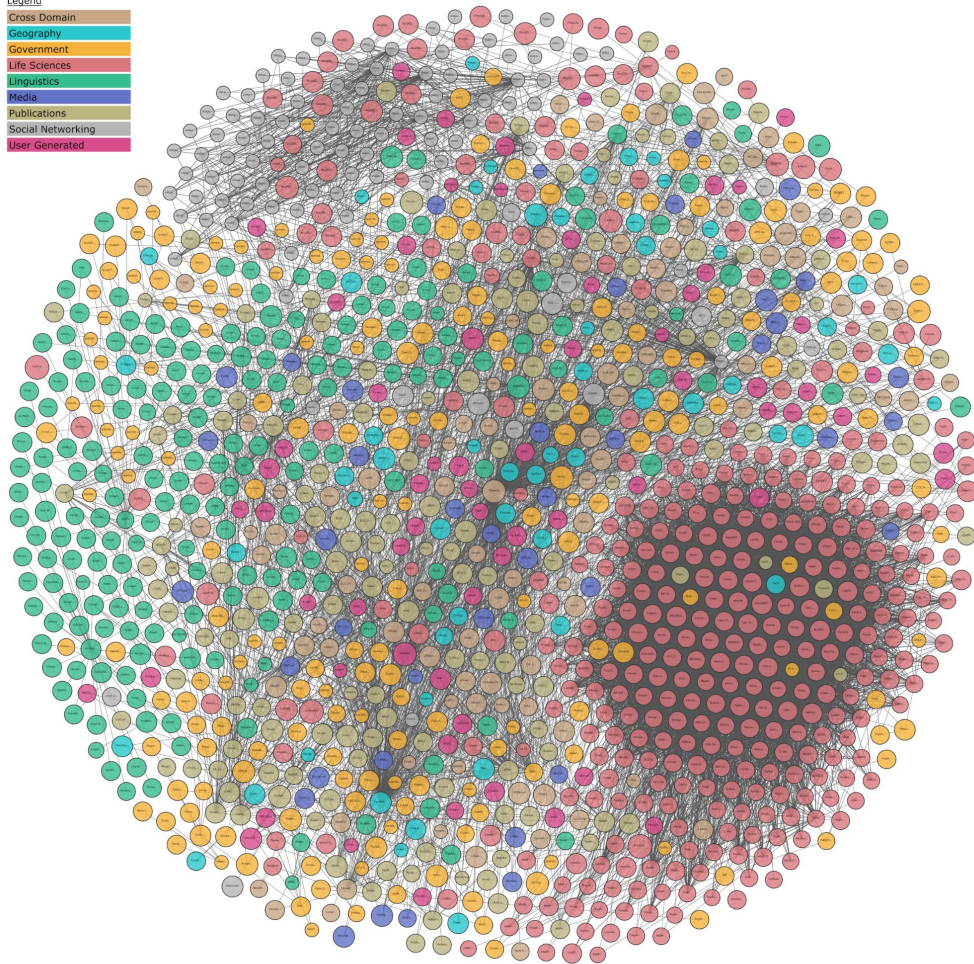
(The result can only be sent back to browser, not saved on the server, see details)

Run Query    Reset

| gene | entrez_id | uniprot_id |
|---|---|---|
| http://dbpedia.org/resource/DsbA | "948353" | "P0AEG4" |
| http://dbpedia.org/resource/Cholinesterase | "590" | "P06276" |
| http://dbpedia.org/resource/Cholinesterase | "590" | "P22303" |
| http://dbpedia.org/resource/Cholinesterase | "43" | "P06276" |
| http://dbpedia.org/resource/Cholinesterase | "43" | "P22303" |
| http://dbpedia.org/resource/Clostridium_perfringens_alpha_toxin | "988262" | |
| http://dbpedia.org/resource/Lymphotoxin | "4049" | "P01374" |
| http://dbpedia.org/resource/Lymphotoxin | "4049" | "Q06643" |
| http://dbpedia.org/resource/Lymphotoxin | "4050" | "P01374" |
| http://dbpedia.org/resource/Lymphotoxin | "4050" | "Q06643" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1457" | "P19784" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1457" | "P67870" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1457" | "P68400" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1460" | "P19784" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1460" | "P67870" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1460" | "P68400" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1459" | "P19784" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1459" | "P67870" |
| http://dbpedia.org/resource/Casein_kinase_2 | "1459" | "P68400" |
| http://dbpedia.org/resource/Collagenase | "4317" | "P03956" |
| http://dbpedia.org/resource/Collagenase | "4317" | "P22894" |
| http://dbpedia.org/resource/Collagenase | "4312" | "P03956" |
| http://dbpedia.org/resource/Collagenase | "4312" | "P22894" |
| http://dbpedia.org/resource/Guanylin | "2980" | "Q02747" |
| http://dbpedia.org/resource/Macrophage_inflammatory_protein | "6348" | "P10147" |
| http://dbpedia.org/resource/Macrophage_inflammatory_protein | "6348" | "P13236" |
| http://dbpedia.org/resource/Macrophage_inflammatory_protein | "6351" | "P10147" |
| http://dbpedia.org/resource/Macrophage_inflammatory_protein | "6351" | "P13236" |

Legend
Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networking
User Generated

The Linked Open Data Cloud from lod-cloud.net

29

Not only DBpedia

31 billion edges in 2011

Named graph http://sparql.uniprot.org/uniprot

| Unique triples | 53332406233 |
| Distinct subjects | 11329433709 |
| Distinct predicates | 122 |
| Distinct objects | 10402001238 |
| Distinct classes | 120 |
| FTP download | |

Overview diagram

Main concepts

53 billion triples in Uniprot KG (Jan. 2024)

# Representing Knowledge Graphs

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

# RDF: formalize statements

Simple sentence "subject - verb - object"

*RAC1 is a human gene.*

Graphical representation



Machine readable syntax (RDF)

`<http://RAC1>    <http://is_a>    <http://Human_Gene> .`

# RDF: syntax to describe (and link) resources

## Definition

(1) an RDF statement represents a **relationship** between two entities: the **subject** and the **object**

(2) the **predicate** represents the nature of their relationship

(3) the relationship is phrased in a **directional** way (from subject to object) and is called in RDF a **property**

(4) RDF statements are called **triples**: they consist of three elements they

(5) Nodes are **URIs** to identify **named entities** on the web or **Literals** to represents text, numbers



`<http://RAC1> <http://is_a> <http://Human_Gene> .`



`<http://snp123> <http://is_a_variant_of> <http://RAC1> .`



`<http://RAC1>  <http://has_amino_acids>  192 .`

# RDF graphs

**Definitions**

(1) A **graph** structure is formed with a set of **nodes** (resources) and **edges** (relationships between resources)

(2) A set of RDF triples is called an RDF graph. RDF is a **directed**, **labeled graph** data format for representing information in the Web.

# Writing RDF graphs with Turtle

## Definitions

(1) One line per triple, each element separated by **space**, each triple endswith a **.**

$$S \quad P \quad O \quad .$$

(2) If two triples describe the same subject, you can reuse it:

$$S \quad P_1 \quad O_1 \quad ;$$
$$P_2 \quad O_2 \quad .$$

(3) If two triples describe the same subject and predicate, you can reuse it:

$$S \quad P \quad O_1 \quad , \quad O_2 \quad .$$



```
@prefix ns: <http://my/namespace/> .

ns:RAC1      rdf:type              ns:Human_gene .
ns:snp123    ns:is_a_variant_of    ns:RAC1 ;
             ns:is_a_risk_factor_of ns:Mitral_stenosis
.
```

# In practice … from text to KG

**Exercise**

From wikipedia : *"The insulin receptor (IR) is a transmembrane receptor that is activated by insulin, IGF-I, IGF-II and belongs to the large class of receptor tyrosine kinase."*

1. Draft a **graphical representation** of the associated knowledge graph.
   - ✓ Identify verbs → predicates
   - ✓ Identify linked entities,
     who is a subject of a relation ?
     who is the object of a relation ?

2. Give the **RDF syntax** for this KG.

# In practice …

*"The insulin receptor (IR) is a transmembrane receptor that is activated by insulin, IGF-I, IGF-II and belongs to the large class of receptor tyrosine kinase."*

# In practice …



```
@prefix ns: <http://my/namespace/> .

ns:insulin   ns:activates   ns:IR .
ns:IGF_I     ns:activates   ns:IR .
ns:IGF_II    ns:activates   ns:IR .

ns:IR     rdf:type              ns:TransmembraneReceptor ;
          Ns:belongs_to_class  ns:ReceptorTyrosineKinase .
```

# Practice ...

**You**

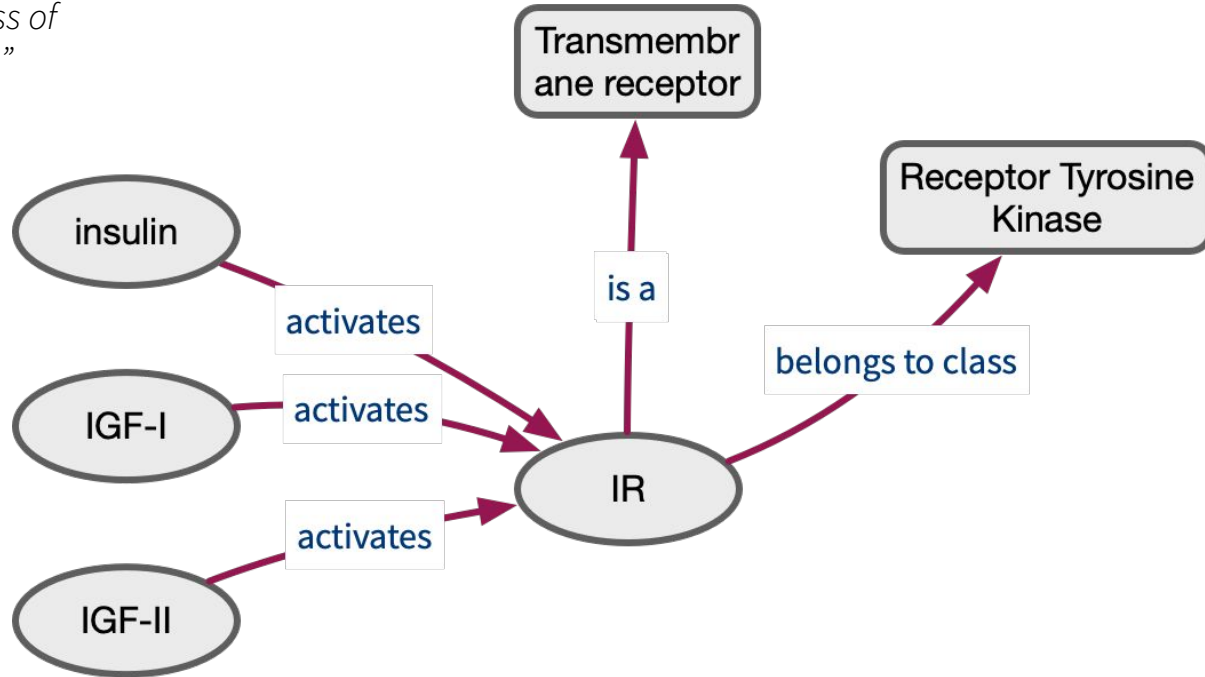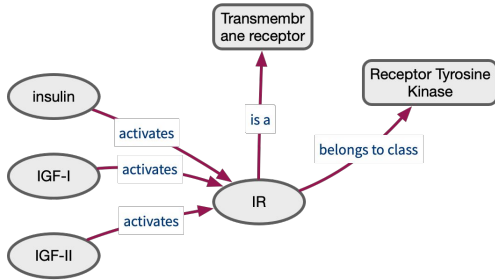Can you translate in RDF (turtle syntax) the following statement : "The insulin receptor (IR) is a transmembrane receptor that is activated by insulin, IGF-I, IGF-II and belongs to the large class of receptor tyrosine kinase"

**ChatGPT**

Sure, here's how you could represent that statement in RDF using Turtle syntax:

```turtle
@prefix ex: <http://example.org/> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix owl: <http://www.w3.org/2002/07/owl#> .

ex:InsulinReceptor rdf:type ex:TransmembraneReceptor ;
    ex:activatedBy ex:Insulin,
                   ex:IGF-I,
                   ex:IGF-II ;
    ex:belongsToClass ex:ReceptorTyrosineKinase .

ex:Insulin a ex:Hormone .
ex:IGF-I a ex:Hormone .
ex:IGF-II a ex:Hormone .
```

# In practice … from KG to text

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix up: <http://purl.uniprot.org/core/> .
@prefix annotation: <http://purl.uniprot.org/annotation/> .

<P06213> rdf:type up:Protein ;
        up:citation citation:2859121 ,
                     citation:2983222 ,
        up:annotation annotation:PRO_0000016687 ,
                      annotation:PRO_0000016689 ,
                      annotation:VAR_015924 .

citation:2859121 rdf:type up:Journal_Citation ;
   up:title "The human insulin receptor cDNA: the structural basis for
hormone-activated transmembrane signalling." ;
   up:author "Ebina Y." , "Ellis L." ;
   skos:exactMatch pubmed:2859121 .

annotation:PRO_0000016687 rdf:type up:Chain_Annotation ;
   rdfs:comment "Insulin receptor subunit alpha" ;
   up:mass 83642 ;
   up:range range:22571007465304878tt28tt758 .
range:22571007465304878tt28tt758 rdf:type faldo:Region ;
   faldo:begin position:22571007465304878tt28 ;
   faldo:end position:22571007465304878tt758 .
```
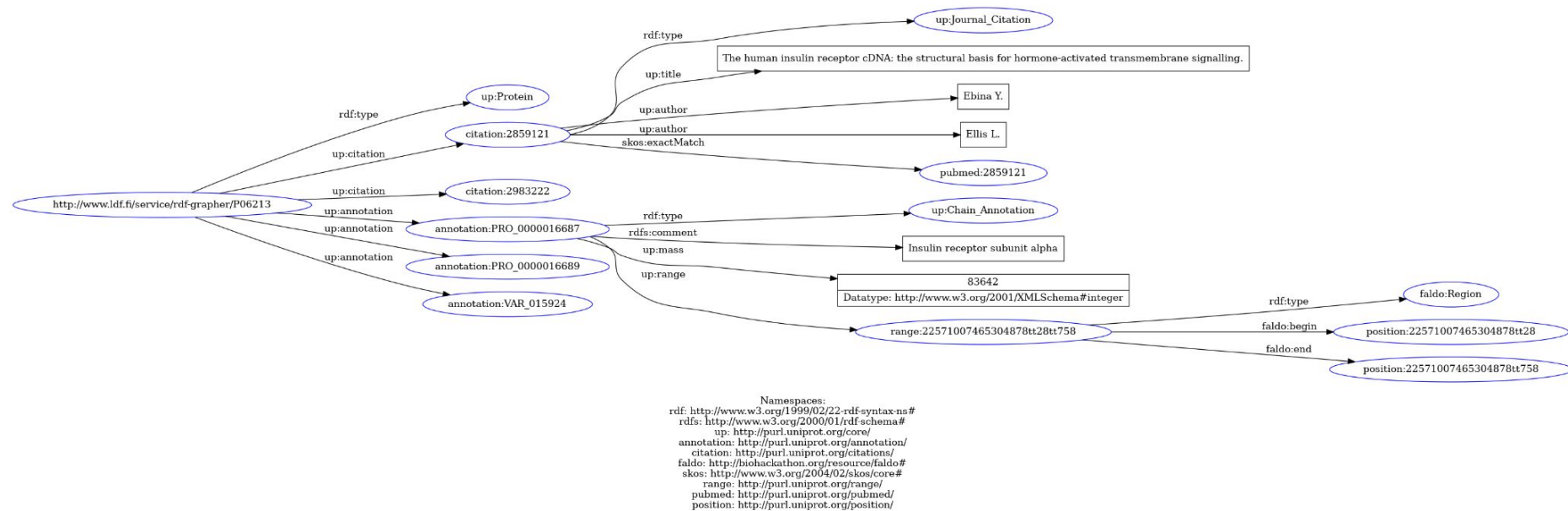
## Exercice

1. Draft the knowledge graph associated to the RDF triples of the P06213 Uniprot entity.
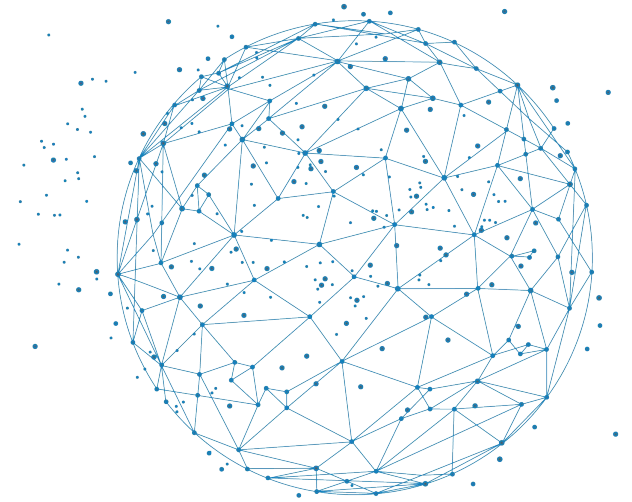
43

# Practice … from KG to text

Namespaces:
rdf: http://www.w3.org/1999/02/22-rdf-syntax-ns#
rdfs: http://www.w3.org/2000/01/rdf-schema#
up: http://purl.uniprot.org/core/
annotation: http://purl.uniprot.org/annotation/
citation: http://purl.uniprot.org/citations/
faldo: http://biohackathon.org/resource/faldo#
skos: http://www.w3.org/2004/02/skos/core#
range: http://purl.uniprot.org/range/
pubmed: http://purl.uniprot.org/pubmed/
position: http://purl.uniprot.org/position/

# Practice … from KG to text

**ChatGPT**

La protéine P06213 est une protéine. Elle est citée dans des journaux ayant pour titre "The human insulin receptor cDNA: the structural basis for hormone-activated transmembrane signaling." Ses auteurs sont Ebina Y. et Ellis L. Cette protéine a des annotations incluant "Sous-unité alpha du récepteur à l'insuline" avec une masse de 83642. Cette annotation correspond à une région spécifique identifiée aux positions 28 à 758.

!! inaccurate positions

!! no link to the pubmed paper

# Querying with graph patterns

# Triple patterns (TP)

SPARQL is the W3C language to query multiple data sources expressed in RDF.

The principle consists in defining a graph pattern to be matched against an RDF graph.

> **Definition**
>
> **Triple Patterns** (TPs) are like RDF triples except that each of the subject, predicate and object may be a **variable**. Variables are prefixed with a **?** .

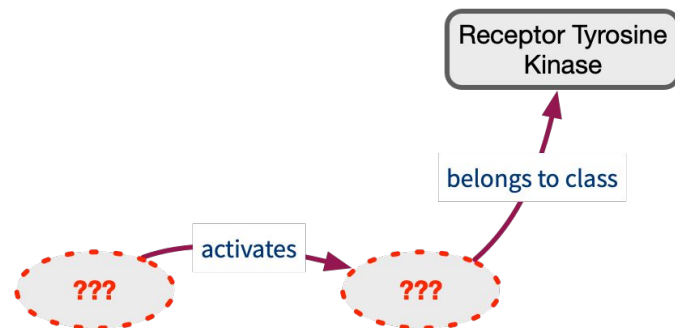*Give me all known activators of IR ?*



```
?x    ns:activates    ns:IR .
```

# Basic graph patterns (BGP)

### Definition

**Basic Graph Patterns** (BGPs) consist in a set of triple patterns to be matched on an RDF graph.

Give me **all** known activators of **any** Receptor Tyrosine Kinase ?

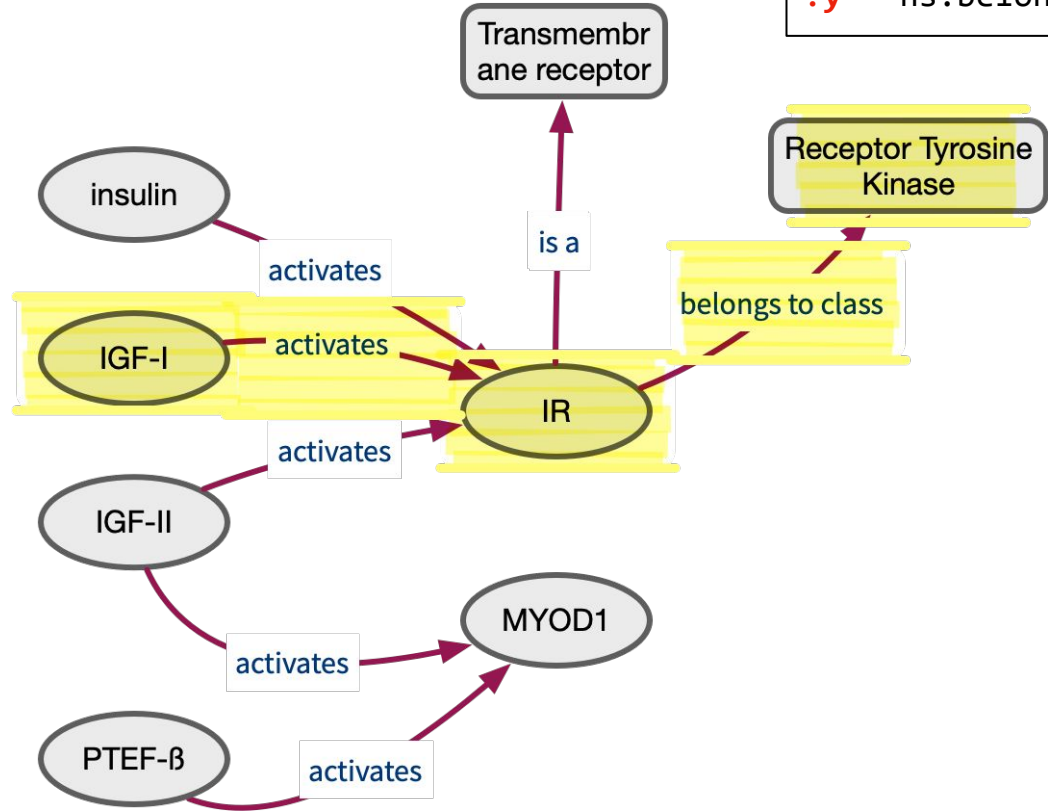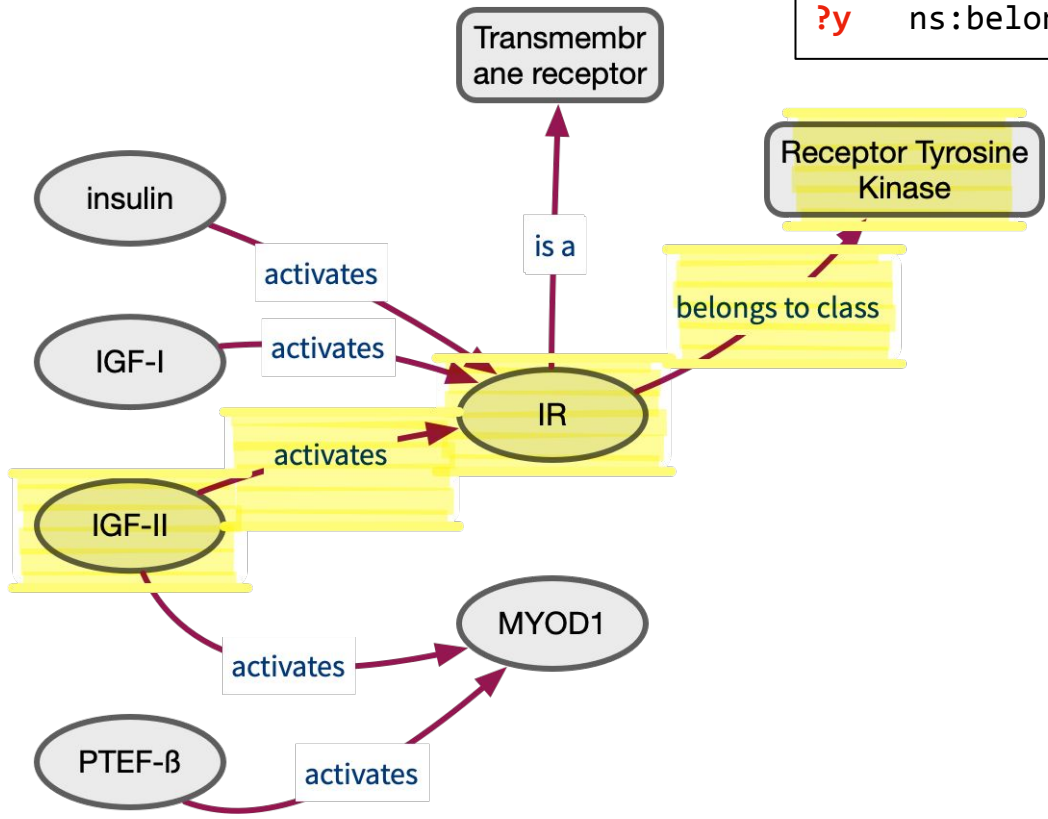→ **all** entities that **activate something** that **belongs to class** "Receptor Tyrosine Kinase"



```
?x   ns:activates    ?y .
?y   ns:belongs_to_class ns:ReceptorTyrosineKinase .
```

```
?x    ns:activates     ?y .
?y    ns:belongs_to_class ns:ReceptorTyrosineKinase .
```



| ?x | ?y |
|---|---|
|  |  |
|  |  |
|  |  |

```
?x    ns:activates      ?y .
?y    ns:belongs_to_class ns:ReceptorTyrosineKinase .
```

| ?x | ?y |
|---|---|
| insulin | IR |
| | |
| | |

```
?x   ns:activates      ?y .
?y   ns:belongs_to_class ns:ReceptorTyrosineKinase .
```

| ?x | ?y |
|---|---|
| IGF-I | IR |
| insulin | IR |
|  |  |

```
?x    ns:activates     ?y .
?y    ns:belongs_to_class ns:ReceptorTyrosineKinase .
```

| ?x | ?y |
|---|---|
| IGF-II | IR |
| IGF-I | IR |
| insulin | IR |

52

Shortcuts
definition

Query clause

BGP

BGP

```
1 ▾ PREFIX rdf:      <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
2   PREFIX rdfs:     <http://www.w3.org/2000/01/rdf-schema#>
3   PREFIX dc:       <http://purl.org/dc/elements/1.1/>
4   PREFIX wp:       <http://vocabularies.wikipathways.org/wp#>
5   PREFIX dcterms:  <http://purl.org/dc/terms/>
6   PREFIX identifiers:<http://identifiers.org/ensembl/>
7   PREFIX atlas:    <http://rdf.ebi.ac.uk/resource/atlas/>
8   PREFIX atlasterms: <http://rdf.ebi.ac.uk/terms/atlas/>
9   PREFIX efo: <http://www.ebi.ac.uk/efo/>
10
11 ▾ SELECT DISTINCT ?wpURL ?pwTitle ?expressionValue ?pvalue where {
12
13 ▾ SERVICE <https://www.ebi.ac.uk/rdf/services/atlas/sparql> {
14       ?factor rdf:type efo:EFO_0000270 .
15       ?value atlasterms:hasFactorValue ?factor .
16       ?value atlasterms:isMeasurementOf ?probe .
17       ?value atlasterms:pValue ?pvalue .
18       ?value rdfs:label ?expressionValue .
19       ?probe atlasterms:dbXref ?dbXref .
20   }
21       ?pwElement dcterms:isPartOf ?pathway .
22       ?pathway dc:title ?pwTitle .
23       ?pathway dc:identifier ?wpURL .
24       ?pwElement wp:bdbEnsembl ?dbXref .
25   }
26   ORDER BY ASC(?pvalue)
```

Query
pattern

modifier

53

# Reasoning with Knowledge graphs

**Handle synonyms** (from PubMed https://pubmed.ncbi.nlm.nih.gov/)

- Look for articles about "vitamin c" in full text search
- Look at the MeSH annotations
- Look for the MeSH term vitamin C and the articles it annotates
- Look for the MeSH term ascorbic acid and the articles it annotates

**Handle taxonomy** (from the MeSH https://www.nlm.nih.gov/mesh/)

- Look for cardiovascular disease
- Select the relevant MeSH term ( https://meshb.nlm.nih.gov/record/ui?ui=D002318 )
- Look at its synonyms and its descendants
- Add it to the search builder
- Search on PubMed

Synonyms and taxonomy are handled transparently

In the GO website (http://geneontology.org/)

- Look for "glucose metabolic process"
- Select "ontology" in the radio box
- Select the relevant GO term (http://amigo.geneontology.org/amigo/term/GO:0006006)
- Select either the "graph view" or the "inferred tree view"
  - Visualise the GO term ancestors
  - Visualize the GO term descendants
- For Homo sapiens, how may proteins, miRNA, etc are annotated by this GO term (or one of its descendants)?

# What is an ontology ?

**Explicit** → a machine does not spontaneously "understand", "infer" or "reason"

**Conceptualization** → a knowledge model aimed at reducing the complexity (generalizing) real facts

**Formal** → reasoning mechanisms must be correct for reliable deductions

**Shared** → domain knowledge result from the consensus of expert communities

## Definitions

*"[...] an explicit specification of a conceptualization"* (Gruber, 1993)

*"[...] a formal specification of a shared conceptualization"* (Borst, 1997)

# Toy example

The **insulin receptor** (**IR**) is a transmembrane receptor that is activated by insulin, IGF-I, IGF-II and belongs to the large class of receptor tyrosine kinase.[5] Metabolically, the insulin receptor plays a key role in the regulation of glucose homeostasis, a functional process that under degenerate conditions may result in a range of clinical manifestations including diabetes and cancer.

Transmembrane receptor

Receptor Tyrosine Kinase

Regulation of glucose homeostasis

Functional process

Clinical manifestation

How these concepts are related together ?

is activated by

results in

results from

How these relations link concepts together ?

Do they allow deductions ?

# RDF-S to represent simple ontologies

**RDF-Schema** aims at providing a simple vocabulary to **organize domain-specific knowledge** through classes (**concepts**) and properties (**relationships**).

## Class VS Instances

Resources may be classified into groups called **classes**. The members of a class are known as **instances** of the class. The rdf:**type** property is used to state that a resource is an instance of a class (« is a » relation ).

## Defining ontologies

- rdf:**type**: to state that a resource is an instance of a class

- owl:**Class** & owl:**Property** to define specific classes or properties

- rdfs:**subClassOf**: to state that all the instances of one class are instances of another

- rdfs:**subPropertyOf**: to state that all resources related by one property are also related by another

- rdfs:**range**: a constraint on the class membership(s) for values of this property

- rdfs:**domain**: a constraint on the class membership(s) for resources having this property

- rdfs:**label**, rdfs:**comment**

# Sample RDF-S vocabulary

```
@prefix etbii: <http://our-namespace#> .
@prefix wikipedia: <https://en.wikipedia.org/wiki/>

etbii:TransmembraneReceptor rdf:type owl:Class ;
        rdfs:subClassOf etbii:Receptor ;
        rdfs:seeAlso wikipedia:Cell_surface_receptor .

etbii:Receptor rdf:type owl:Class ;
            rdfs:subClassOf etbii:Protein .

etbii:Protein rdf:type owl:Class .

etbii:isActivatedBy rdf:type owl:Property .
        rdfs:domain etbii:Protein ;
        rdfs:range  etbii:Protein .
```

# Reasoning with RDFS-Entailments

**Inference rules** to produce new logical facts, or to check for logical soundness (satisfiability)

**Deduce** the multiple **types** of an entity based on **class hierarchies**

**Deduce** the **types** of entities exploiting the definition of **relations**

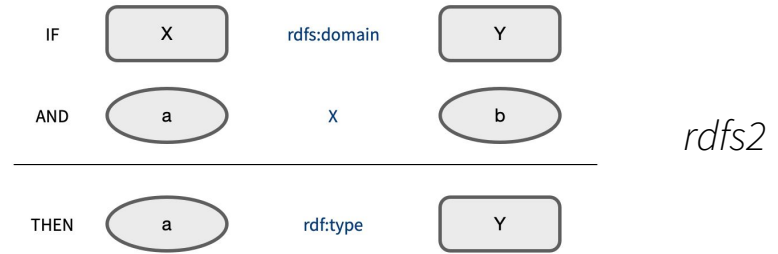… more possibilities with OWL (Web Ontology Language) and Description Logics (DL)

*From the specification of a "link", how to infer the type of a linked entity ?*

**RDFS entailment patterns.**

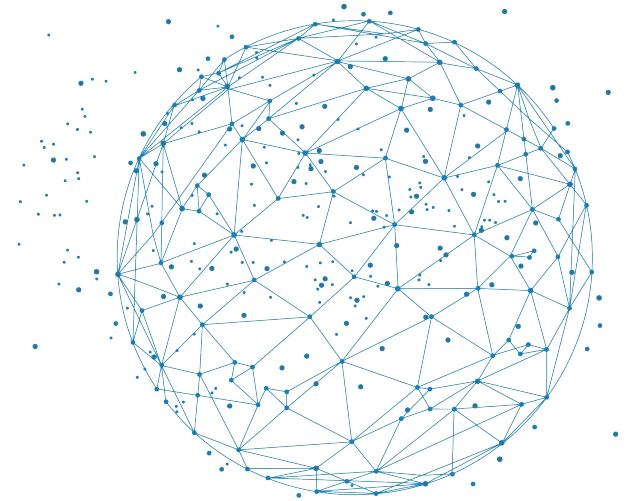| | If S contains: | then S RDFS entails recognizing D: |
|---|---|---|
| *rdfs1* | any IRI aaa in D | aaa rdf:type rdfs:Datatype . |
| *rdfs2* | aaa rdfs:domain xxx .<br>yyy aaa zzz . | yyy rdf:type xxx . |
| *rdfs3* | aaa rdfs:range xxx .<br>yyy aaa zzz . | zzz rdf:type xxx . |
| *rdfs4a* | xxx aaa yyy . | xxx rdf:type rdfs:Resource . |
| *rdfs4b* | xxx aaa yyy. | yyy rdf:type rdfs:Resource . |
| *rdfs5* | xxx rdfs:subPropertyOf yyy .<br>yyy rdfs:subPropertyOf zzz . | xxx rdfs:subPropertyOf zzz . |
| *rdfs6* | xxx rdf:type rdf:Property . | xxx rdfs:subPropertyOf xxx . |
| *rdfs7* | aaa rdfs:subPropertyOf bbb .<br>xxx aaa yyy . | xxx bbb yyy . |
| *rdfs8* | xxx rdf:type rdfs:Class . | xxx rdfs:subClassOf rdfs:Resource . |
| *rdfs9* | xxx rdfs:subClassOf yyy .<br>zzz rdf:type xxx . | zzz rdf:type yyy . |
| *rdfs10* | xxx rdf:type rdfs:Class . | xxx rdfs:subClassOf xxx . |
| *rdfs11* | xxx rdfs:subClassOf yyy .<br>yyy rdfs:subClassOf zzz . | xxx rdfs:subClassOf zzz . |
| *rdfs12* | xxx rdf:type rdfs:ContainerMembershipProperty . | xxx rdfs:subPropertyOf rdfs:member . |
| *rdfs13* | xxx rdf:type rdfs:Datatype . | xxx rdfs:subClassOf rdfs:Literal . |



*rdfs2*

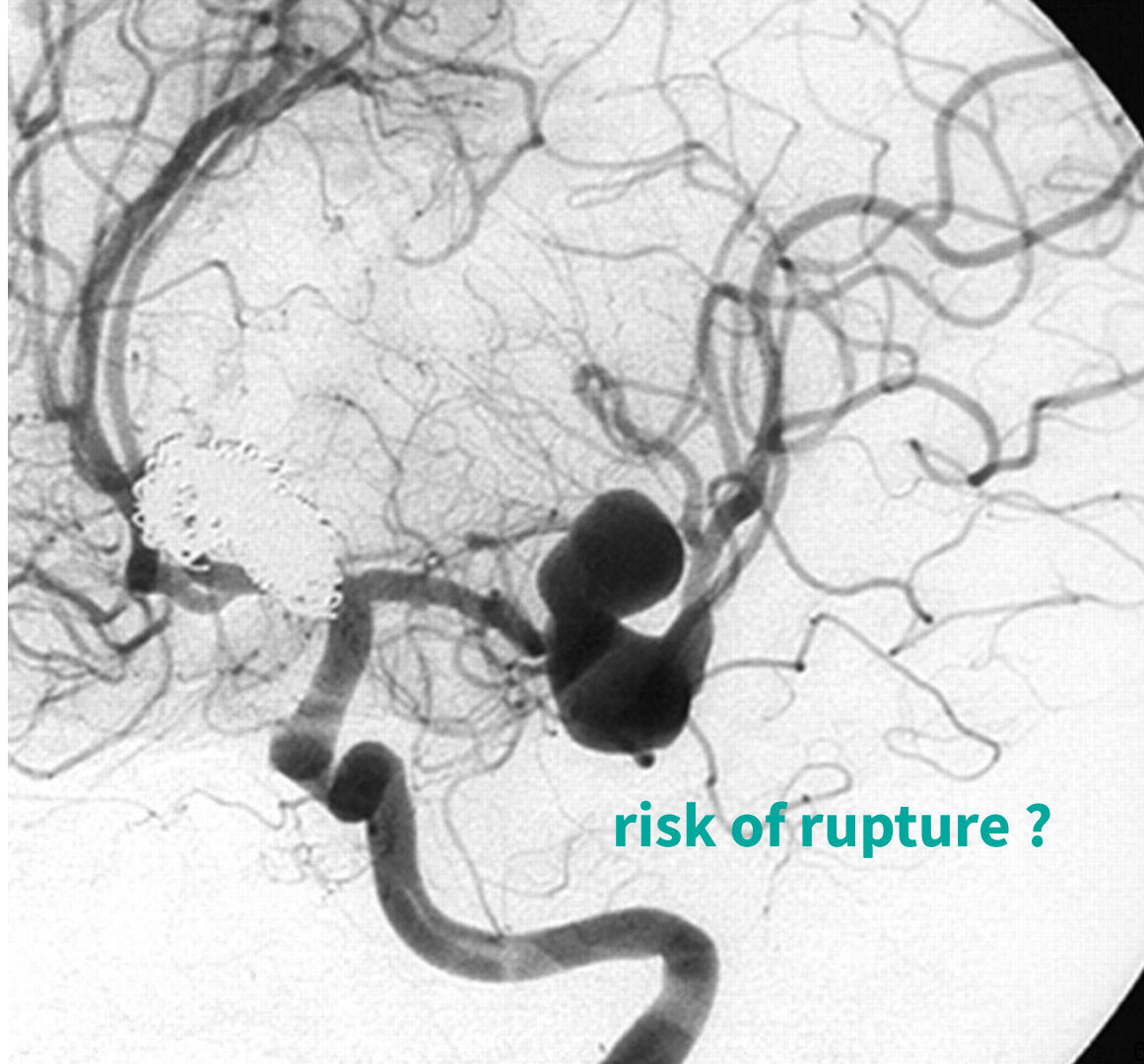*rdfs2, rdfs3*

a = IR
b = insulin
X = is activated by
Y = Protein

# Application with multi-scale data

**ICAN cohort : 34 univ. hospitals / 3400 subjects**

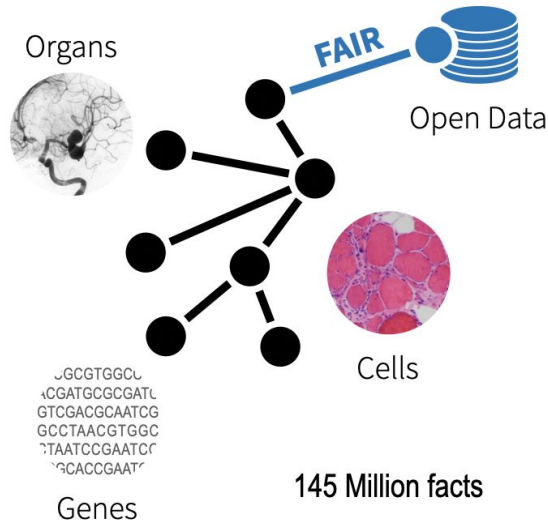- 3000 MRIs
- 800 whole genomes (under processing)

**risk of rupture ?**

# INEX-MED

**Intracranial aneurysms** (3.2 % world population)
**Congenital myopathies** (rare diseases)

Organs

FAIR

Open Data

Cells

GCGTGGC
CGATGCGCAT
GTCGACGCAATCG
GCCTAACGTGGC
TAATCCGAATC
GCACCGAAT

Genes

145 Million facts

**Multi-scale « Knowledge Graph »**

**Select
Combine
Reuse
(Share)**

Ontologies
SPARQL queries

**Predict
Explain**

Statistics
Reasoning
Machine learning

@Nantes

@Strasbourg

Association between imaging phenotypes - « omics » signatures ?
Patients with high/low aneurysm rupture risk ?

# Quantifications of cerebral artery bifurcations

*Anass Nouri, Florent Autrusseau*



- distances ?
- angles ?
- radius ?
- tortuosity ?

Anass Nouri

# Quantifications of cerebral artery bifurcations

*Anass Nouri, Florent Autrusseau*



interpretating these values ?
comparing them ?
sharing ?

$d_1, d_2$ : **distances** between **bifurcations**

$r_1, r_2, r_3$ : **radius** of vessels

$\theta_1, \theta_2$ : **angles** of bifurcations



Anass Nouri

# BrainHack project

« From MS excel sheets to semantic bioimaging markers : representing and querying cerebral vascular measures with the QIBO ontology »

```
SELECT * WHERE {
  # Search variants
  ?subject ican:has_dna_sample ?sample .
  ?sample ican:is_heterozygotous ?variant .
  ?variant ican:has_chromosome "19" ;
           ican:has_position ?position ;
           ican:has_alt_allele ?alt_allele ;
           ican:has_ref_allele ?ref_allele .
  FILTER (?position > 10092337 && ?position < 10106407)
  # dbSNP identifier
  OPTIONAL { ?variant ican:has_dbsnp_id ?dbSNPid . }
  # First aneurysm rupture: age and localisation
  OPTIONAL { ?subject ican:age_first_rupture ?ageRupt . }
  OPTIONAL { ?subject ican:loc_first_rupture ?locRupt . }
}
```

# KG Application in ML

# Graph embeddings approach

from a graph representation ...



to real vector representation



embedding algorithm

from https://gearons.org

**Some GNN Tasks**

Disease-Gene association
Protein Function prediction
PPI prediction
Drug-Target (Drug) Interaction prediction
Drug response …

Zhang X-M, Liang L, Liu L and Tang M-J (2021) Graph Neural Networks and Their Current Applications in Bioinformatics. Front. Genet. 12:690049. doi: 10.3389/fgene.2021.690049

**Example in practice**



| Input (KG facts) | Encoder (GNNs) | Decoder (KGEs) | Output (s,r,o embeddings) |

Zhang X-M, Liang L, Liu L and Tang M-J (2021) Graph Neural Networks and Their Current Applications in Bioinformatics. *Front. Genet.* 12:690049. https://doi.org/10.3389/fgene.2021.690049



Using graph-based model to identify cell specific synthetic lethal effects
https://doi.org/10.1016/j.csbj.2023.10.011

# Embeddings methods on knowledge graphs give different results



**KG embeddings models**

Uniprot Embeddings

RDF2VEC

DistMult

MultiKE

ComplEx

RDGCN

i-Align

Alshahrani M, Thafar MA, Essack M. 2021. Application and evaluation of knowledge graph embeddings in biomedical data. *PeerJ Computer Science* 7:e341 https://doi.org/10.7717/peerj-cs.341

**Knowledge Graphs (KGs)**

**Cons:**
- Implicit Knowledge
- Hallucination
- Indecisiveness
- Black-box
- Lacking Domain-specific/New Knowledge

**Pros:**
- Structural Knowledge
- Accuracy
- Decisiveness
- Interpretability
- Domain-specific Knowledge
- Evolving Knowledge

**Pros:**
- General Knowledge
- Language Processing
- Generalizability

**Cons:**
- Incompleteness
- Lacking Language Understanding
- Unseen Facts

**Large Language Models (LLMs)**

**RAG Application enriched with a Knowledge Graph**



Retrieval Augmented Generation and Knowledge Graphs

https://gradientflow.com/boosting-llms-with-external-knowledge-the-case-for-knowledge-graphs/

**=> Boosting LLMs with External Knowledge: The Case for Knowledge Graphs**

Pan, S., Luo, L., Wang, Y., Chen, C., Wang, J., & Wu, X. (2023). Unifying Large Language Models and Knowledge Graphs: A Roadmap. *ArXiv, abs/2306.08302.*

79

# Some Life Science application of LLMs + KG



KG-RAG. Soman et al. https://arxiv.org/pdf/2311.17330.pdf

Overview of the SPIRES approach.
SPIRES is available as part of the open source OntoGPT
package: https://github.com/ monarch-initiative/ontogpt

KG to LLM
- Document / explain generated text with "facts" coming from knowledge graphs : e.g. Explicability https://arxiv.org/pdf/2309.01029.pdf https://arxiv.org/abs/2311.09188
- Explain data with complex structure : e.g. map data with ontologies, ontology alignment
- Answer domain-specific questions : e.g. guided fine-tuning

LLM to KG
- Natural language interface to generate SPARQL queries : e.g. SPARQL Generation: an analysis on fine-tuning OpenLLaMA for Question Answering over a Life Science Knowledge Graph. Ana Claudia Sima et al. SWAT4HCLS 2024 *https://arxiv.org/abs/2402.04627*
- Augment knowledge with synthetic data for better prediction in graph embedding approaches e.g. https://arxiv.org/abs/2203.13965



a. KG-enhanced LLMs     b. LLM-augmented KGs     c. Synergized LLMs + KGs

Pan, S., Luo, L., Wang, Y., Chen, C., Wang, J., & Wu, X. (2023). Unifying Large Language Models and Knowledge Graphs: A Roadmap. *ArXiv, abs/2306.08302.*

## Phenomics Assistant: An Interface for LLM-based Biomedical Knowledge Graph Exploration

Shawn T O'Neil, Kevin Schaper, Glass Elsarboukh, Justin T Reese, Sierra A T Moxon, Nomi L Harris, Monica C Munoz-Torres, Peter N Robinson, Melissa A Haendel, Christopher J Mungall

**doi:** https://doi.org/10.1101/2024.01.31.578275

### BioMedGPT: Open Multimodal Generative Pre-trained Transformer for BioMedicine
Yizhen Luo et al. https://arxiv.org/abs/2308.09442

### Structured Prompt Interrogation and Recursive Extraction of Semantics (SPIRES): a method for populating knowledge bases using zero-shot learning 🔓

J Harry Caufield ✉, Harshad Hegde, Vincent Emonet, Nomi L Harris, Marcin P Joachimiak, Nicolas Matentzoglu, HyeongSik Kim, Sierra Moxon, Justin T Reese, Melissa A Haendel … Show more

### Biomedical knowledge graph-enhanced prompt generation for large language models
Soman et al. https://arxiv.org/pdf/2311.17330.pdf

### Biologically informed variational autoencoders allow predictive modeling of genetic and drug−induced perturbations

Daria Doncevic [1], Carl Herrmann [1]

Affiliations + expand
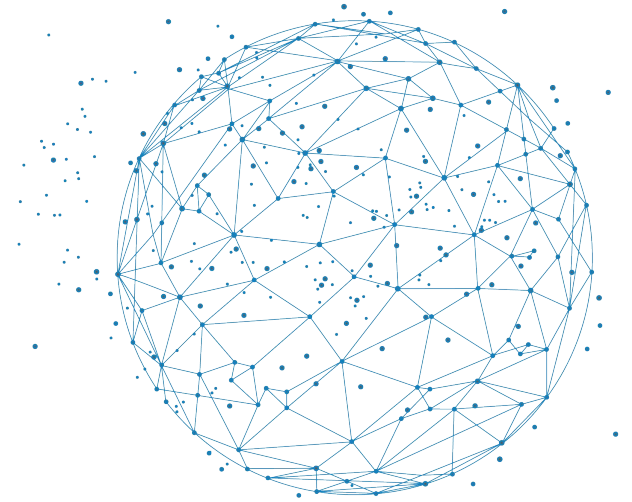PMID: 37326971    PMCID: PMC10301695    DOI: 10.1093/bioinformatics/btad387

### Dynamic Retrieval Augmented Generation of Ontologies using Artificial Intelligence (DRAGON-AI)
Sabrina Toro *et al*
https://arxiv.org/abs/2312.10904

# Keep in mind

*Complex data analyses require fine-grained, explicit descriptions*

- Annotate your data with **RDF** to assemble **knowledge graphs** (KGs)

- Support future **integration** by referring to other Knowledge Graphs: **URIs**

- Formalize domain knowledge with **ontologies**: **RDFS**, **OWL**

- Mine (multiple) KGs with **graph patterns**: (federated) **SPARQL** queries

*Advantages*

- Graphs for humans & machines

- Semantic heterogeneity

- Established web technologies

- Technological framework
  for F - * - I - R principles

- De-centralized: many query-able datasets published on the web (Linked Data Cloud)

*Disadvantages*

- Semantic heterogeneity (many Life Science ontologies …)

- Reliability of external SPARQL endpoints

- No graphical user interfaces… learn RDF syntaxes and SPARQL query language

- Scalable querying is still a hot research topic …

# Going further

- Bob DuCharme
  - What is RDF?
    http://www.bobdc.com/blog/whatisrdf/
  - What is RDFS?
    http://www.bobdc.com/blog/whatisrdfs/
  - SPARQL in 11 minutes
    https://www.youtube.com/watch?v=FvGndkpa4K0
  - Learning SPARQL, 2nd ed. O'Reilly

- https://www.w3.org/TR/rdf11-primer/

- https://www.w3.org/TR/sparql11-query/

- https://www.slideshare.net/LeeFeigenbaum/sparql-cheat-sheet

- http://www.wikipathways.org/index.php/Help:WikiPathways_Sparql_queries

- https://www.fun-mooc.fr/fr/cours/web-semantique-et-web-de-donnees/