

# Lecture 8

## Artificial Intelligence: Knowledge Graphs

Part I: Intelligent Agents, Resource Description Framework (RDF)

COMP 6721, Fall 2023

# Outline

- ① Introduction
- ② The Resource Description Framework (RDF)
- ③ Example: Data Integration with Knowledge Graphs
- ④ Notes and Further Reading

## Slides Credit

- Includes slides from Jay Pujara & Sameer Singh, *Mining Knowledge Graphs from Text*, <https://kgtutorial.github.io/>
- Includes slides by Ivan Herman, W3C [Her]
- Includes slides from Hoifung Poon, Chris Quirk & Scott Wen-Tau Yih, *Machine Reading for Precision Medicine*,  
[https://www.microsoft.com/en-us/research/uploads/prod/2018/01/1802\\_aaai-tutorial\\_precision-med.pdf](https://www.microsoft.com/en-us/research/uploads/prod/2018/01/1802_aaai-tutorial_precision-med.pdf)

# Outline

## 1 Introduction

Motivation

Why Knowledge Graphs?

Example: AI for Precision Medicine

History

## 2 The Resource Description Framework (RDF)

## 3 Example: Data Integration with Knowledge Graphs

## 4 Notes and Further Reading

# Eliza

The screenshot shows a window titled '\*doctor\*' with a menu bar containing File, Edit, Options, Buffers, Tools, Text, and Help. The main area displays a conversation:

I am the psychotherapist.  
Please, describe your problems. Each time you are finished talking, type RET twice.

What is Concordia University?

Why do you say what is concordia university?

Just tell me what it is!

Is it because of your plans that you say just tell you what it is?

[A small black square icon is located near the bottom left of the window.]

# Google Assistant



What is Concordia University?



**Concordia University** (French: Université **Concordia**; commonly referred to as **Concordia**) is a public comprehensive research **university** located in Montreal, Quebec, Canada. ... **Concordia** is a non-sectarian and coeducational institution, with more than 215,000 alumni worldwide.

Athletics brand: Concordia Stingers

Subsidiary or constituent schools: Loyola Campus, Faculty of Arts

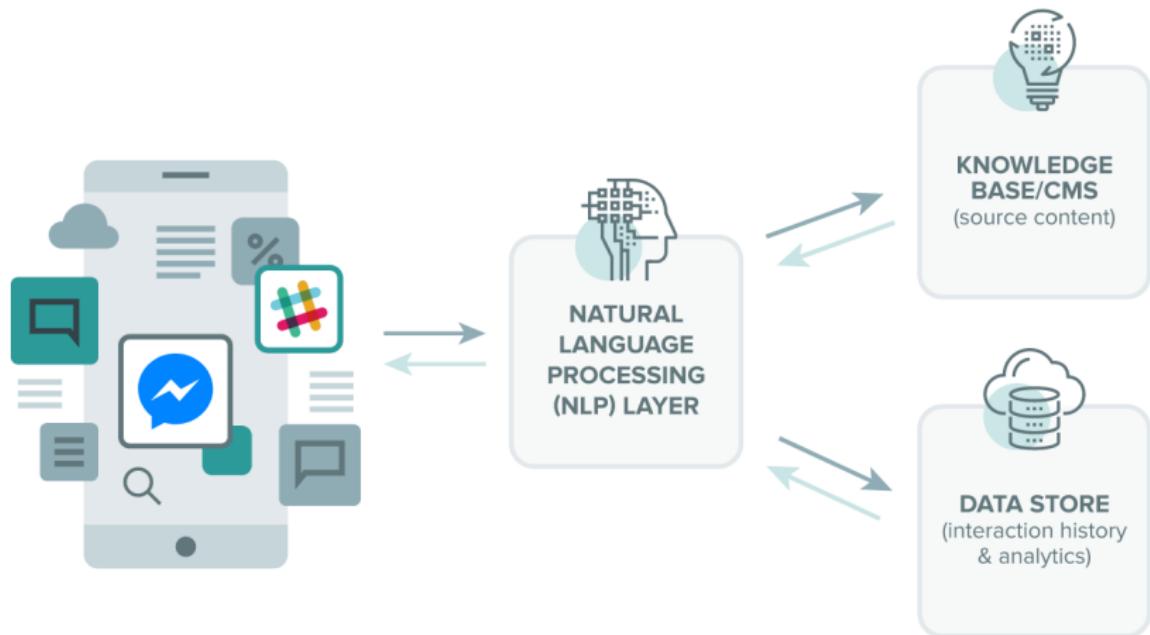
Date founded: August 24, 1974

Geographic scope: Canada

# IBM Watson

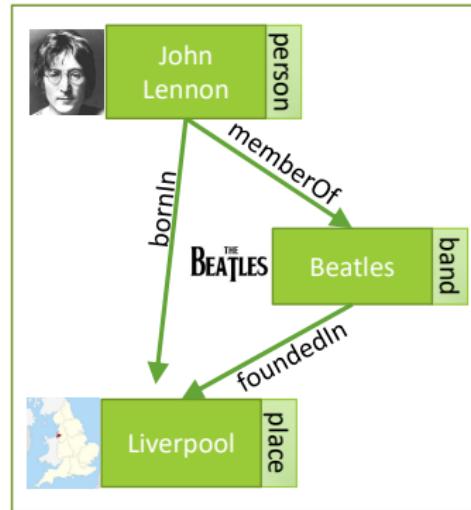


# Generic Intelligent Assistant Architecture



# Example knowledge graph

- Knowledge in graph form!
- Captures entities, attributes, and relationships
- Nodes are entities
- Nodes are labeled with attributes (e.g., types)
- Typed edges between two nodes capture a relationship between entities



# Why knowledge graphs?

---

- Humans:
  - Combat information overload
  - Explore via intuitive structure
  - Tool for supporting knowledge-driven tasks
- AIs:
  - Key ingredient for many AI tasks
  - Bridge from data to human semantics
  - Use decades of work on graph analysis

# Applications 1: QA/Agents



A screenshot of a web search results page for "who is playing in this year's super bowl". The search bar shows the query. Below it, a navigation bar has "All" selected, followed by News, Shopping, Videos, Maps, More, Settings, and Tools. The main content area shows "About 4,350,000 results (0.46 seconds)". A large section titled "Super Bowl LII" displays the NFL game information: "NFL - Today, 3:30 PM", "Philadelphia Eagles" vs "New England Patriots" at the "Super Bowl". It includes a "Game preview" link and a "Watch on NBC" link. A note at the bottom says "All times are in Pacific Time".

# Applications 2: Decision Support

IBM Watson Knowledge Studio

View Details Attribute View View Guidelines Completed (0) Close Alpha... 14pt 1

Entity Mention

Type Subtype Role

- a ACCIDENT\_CAUSE
- e ACCIDENT\_OUTCOME
- y CONDITION
- i IMPACT
- m MANUFACTURER
- Model MODEL
- y MODEL\_YEAR
- r PART\_OF\_CAR
- p PERSON
- s STRUCTURE
- v VEHICLE

2004-49-168A.txt

1 **V1**, a 1996 **Toyota** Camry, was traveling southbound in the second **lane** of a four-lane divided (seven **lanes** overall, divided by raised median), concrete **roadway**, approaching an **intersection**.

2 **V2**, a 2004 Mercedes S430, was northbound in the fourth **lane** of a four-lane, divided (seven **lanes** overall, divided by raised median), concrete **roadway**, about to turn left into westbound traffic at the same **intersection**.

3 As both **vehicles** entered the **intersection**, the **front** of **V1** impacted the **front** of **V2**.

4 **V1** rotated clockwise as **V2** rotated counter-clockwise, and the left side of **V1** impacted the right side of **V2** in a sideslip configuration.

5 Both **vehicles** moved southwest to final rest.

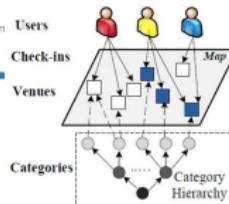
6 Both **vehicles** were towed due to damage.

7 The unrestrained **driver** of **V1** was hospitalized with foot and rib fractures as well as a liver laceration.

8 The restrained **driver** of **V2** was treated and released with minor abrasion and contusion as well as a finger fracture.

9 The restrained **male** right passenger in **V2** was pronounced brain dead two days later from brain injuries.

10 **V1** also suffered with reinforced knee fractures, which developed



# Applications 3: Fueling Discovery

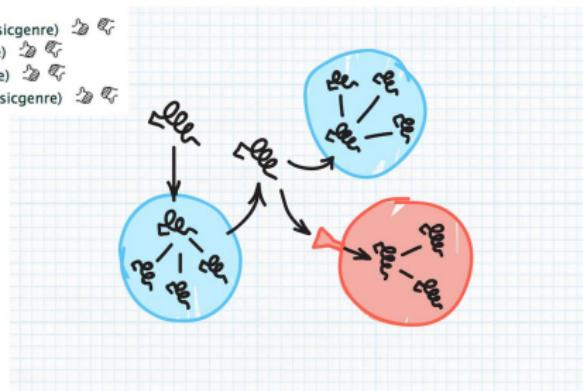
**beatles (musicartist)**

literal strings: BEATLES, Beatles, beatles

**Help NELL Learn!**

NELL wants to know if these be  
If they are or ever were, click thumbs-up. Or

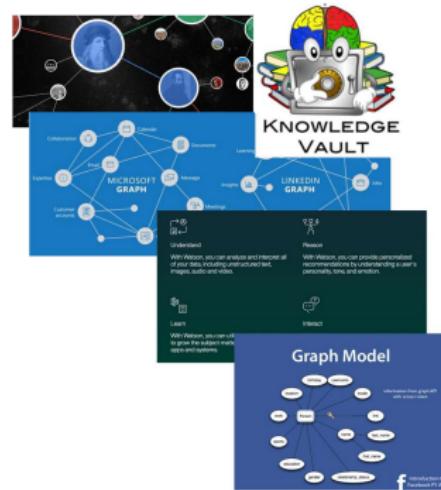
- beatles is a musical artist  
- beatles is a musician in the genre classic pop (musicgenre)  
- beatles is a musician in the genre pop (musicgenre)  
- beatles is a musician in the genre rock (musicgenre)  
- beatles is a musician in the genre classic rock (musicgenre)  



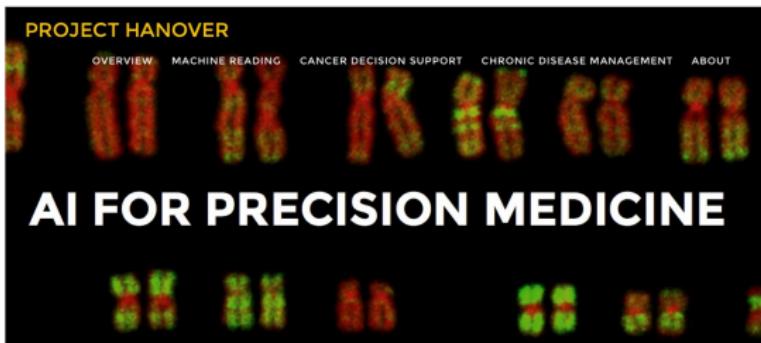
# Knowledge Graphs & Industry

---

- Google Knowledge Graph
  - Google Knowledge Vault
- Amazon Product Graph
- Facebook Graph API
- IBM Watson
- Microsoft Satori
  - Project Hanover/Literome
- LinkedIn Knowledge Graph
- Yandex Object Answer
- Diffbot, GraphIQ, Maana, ParseHub, Reactor Labs, SpazioDati



# Interesting application of Knowledge Graphs

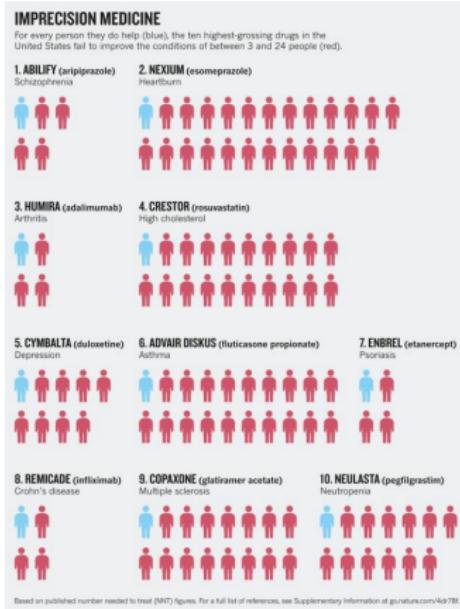


Microsoft<sup>®</sup>  
**Research**

## Chronic disease management:

develop AI technology for predictive and preventive personalized medicine to reduce the national healthcare expenditure on chronic diseases  
(90% of total cost)

# Medicine Today Is Imprecise



Top 20 drugs  
80% non-responders

Wasted  
1/3 health spending  
\$1 Trillion / year

## Example: Tumor Board KB Curation

The deletion mutation on exon-19 of EGFR gene was present in 16 patients, while the L858E point mutation on exon-21 was noted in 10.

All patients were treated with gefitinib and showed a partial response.



Gefitinib can treat tumors w. EGFR-L858E mutation

# PubMed

27 million abstracts

Two new abstracts every minute

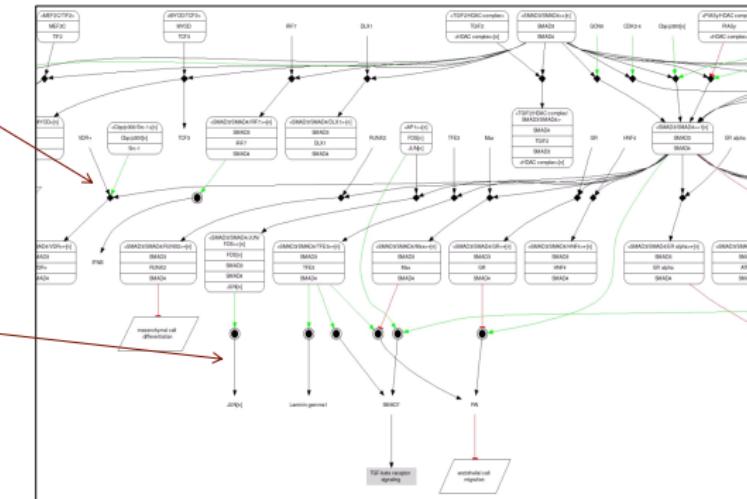
Adds over one million every year



# Machine Reading

PMID: 123  
...  
VDR+ binds to SMAD3 to form  
...  
...

PMID: 456  
...  
JUN expression  
is induced by  
SMAD3/4  
...  
...



# Complex Semantics

Involvement of p70(S6)-kinase activation in IL-10 up-regulation in human monocytes by gp41 envelope protein of human immunodeficiency virus type 1 ...

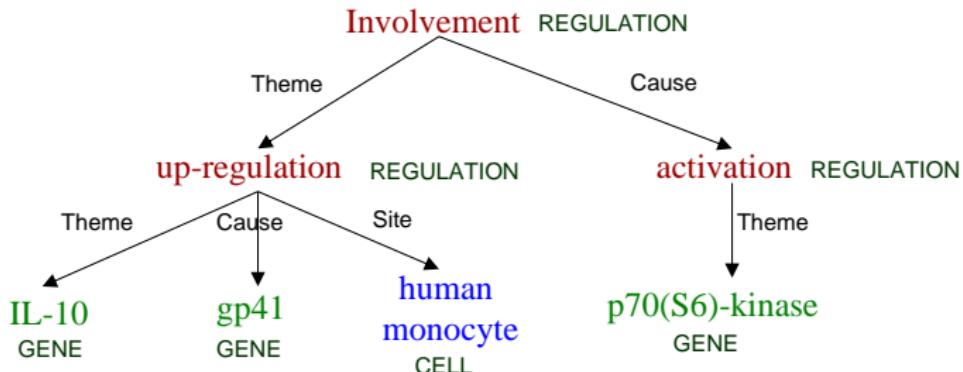
# Complex Semantics

Involvement of p70(S6)-kinase activation in IL-10 up-regulation in human monocytes by gp41 envelope protein of human immunodeficiency virus type 1 ...

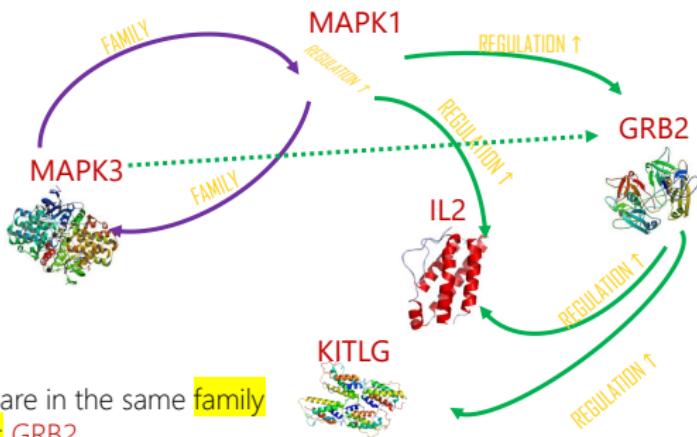
IL-10                  gp41                  human  
GENE                  GENE                  monocyte  
    CELL  
    p70(S6)-kinase  
    GENE

# Complex Semantics

Involvement of p70(S6)-kinase activation in IL-10 up-regulation in human monocytes by gp41 envelope protein of human immunodeficiency virus type 1 ...



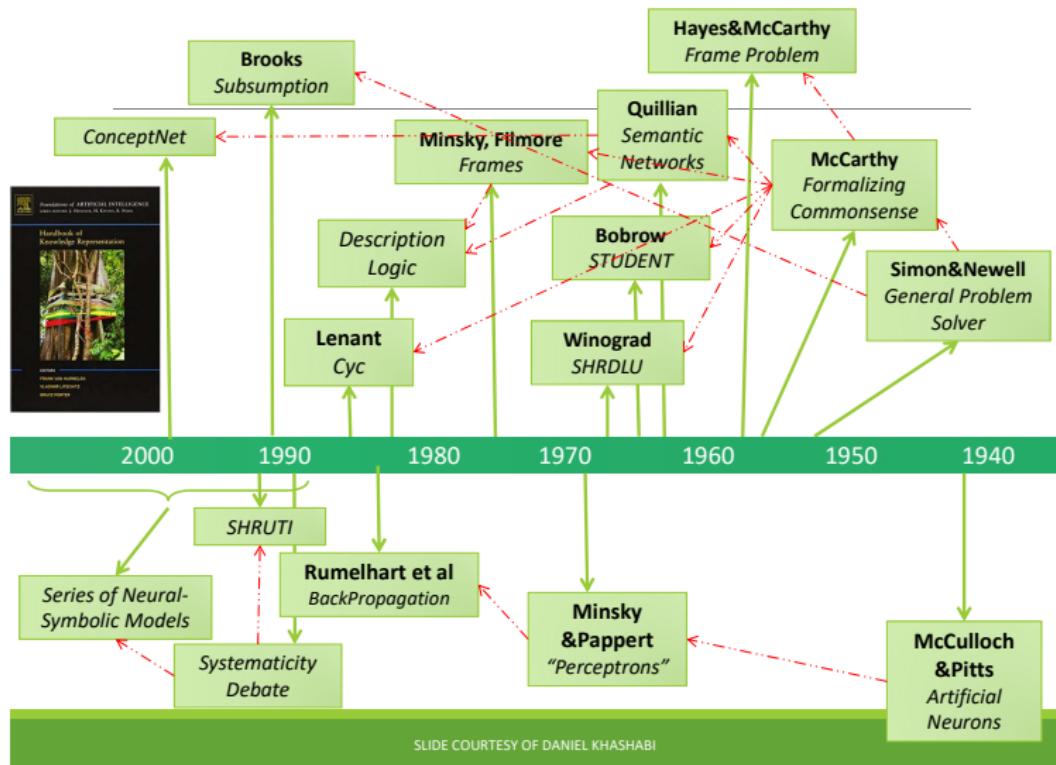
# Genomics Knowledge Base (Network)



MAPK3 and MAPK1 are in the same family  
MAPK1 up-regulates GRB2

Likely that MAPK3 up-regulates GRB2

# History of Knowledge Representation (KR)



# Today

From 1950–2020...

- Concepts have been around for a long time (Semantic Networks, Frames, Description Logic, ...)

# Today

From 1950–2020...

- Concepts have been around for a long time (Semantic Networks, Frames, Description Logic, ...)

1980s/90s

- AI/IS systems suffer from the *Knowledge Acquisition Bottleneck*
- One of the reasons for the *AI Winter* at that time

# Today

From 1950–2020...

- Concepts have been around for a long time (Semantic Networks, Frames, Description Logic, ...)

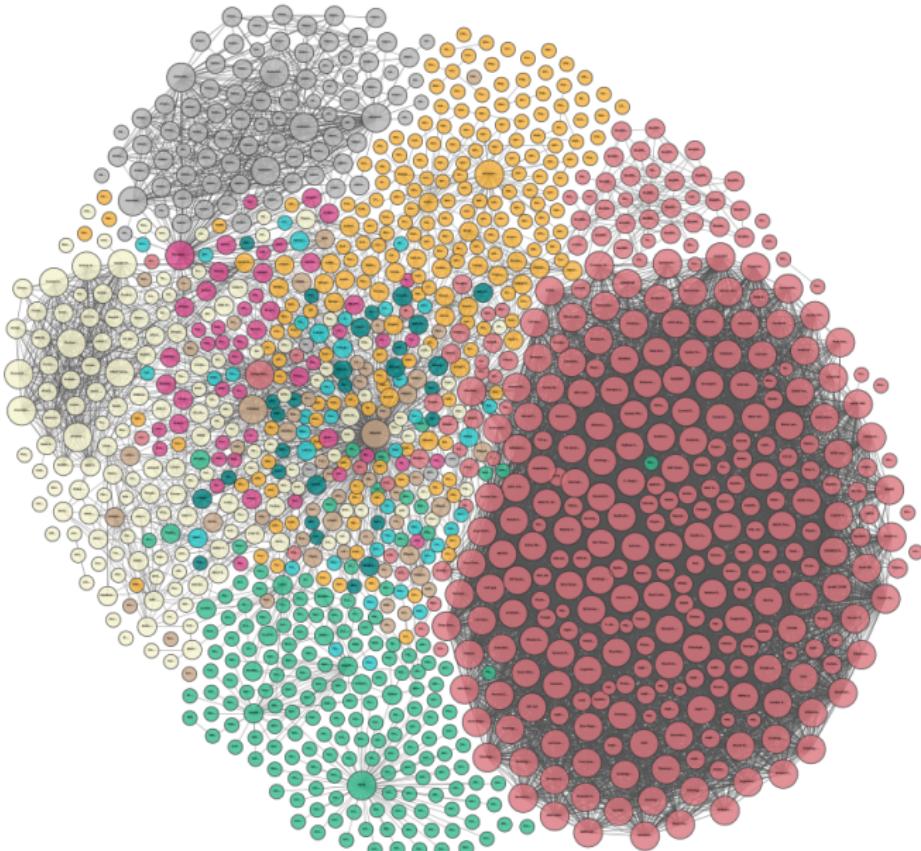
1980s/90s

- AI/IS systems suffer from the *Knowledge Acquisition Bottleneck*
- One of the reasons for the *AI Winter* at that time

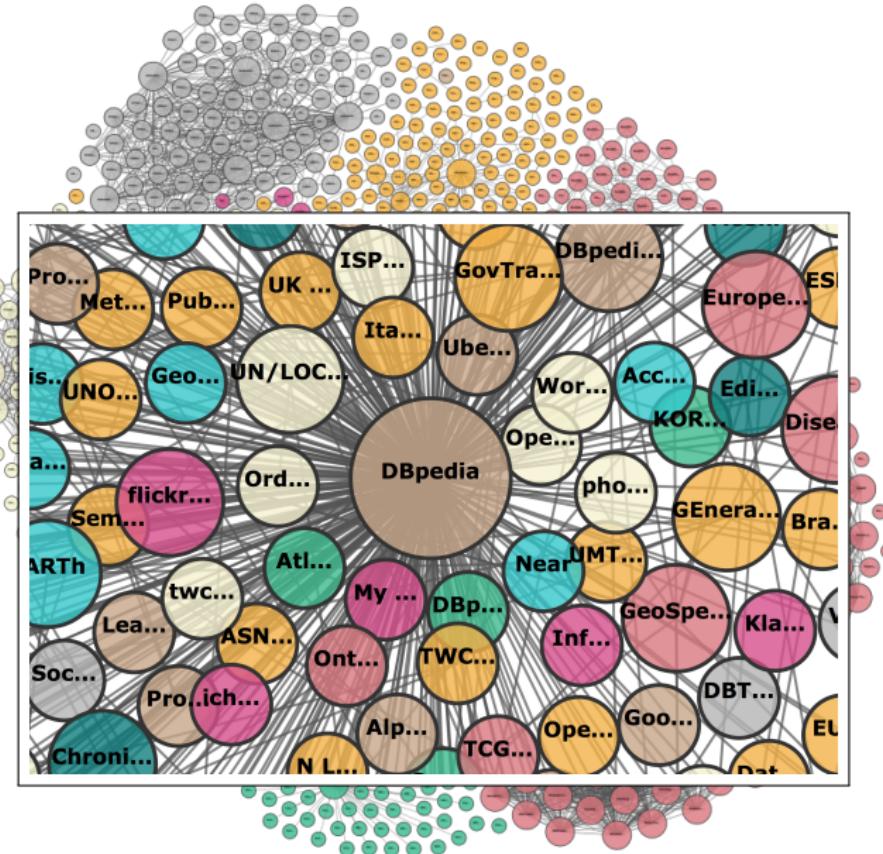
Technology

- Open standards, based on W3C recommendations, e.g., [RDF](#)
- Proprietary products, e.g., [Neo4J](#) or [Oracle Spatial and Graph](#)
- We now have substantial [knowledge bases](#) available, both proprietary  
(e.g., Facebook Graph Search, Google Knowledge Graph) and open access (e.g., Wikidata, DBpedia, YAGO)

# The Linked Open Data Cloud



# The Linked Open Data Cloud



# TBL at TED on “The Next Web” (2009)



Tim Berners-Lee: The next Web of open, linked data

[https://www.youtube.com/watch?v=OM6XIIcm\\_qo](https://www.youtube.com/watch?v=OM6XIIcm_qo)

# Outline

## 1 Introduction

## 2 The Resource Description Framework (RDF)

Introduction

RDF Triples

Literals

Namespaces

Serialization

Programming

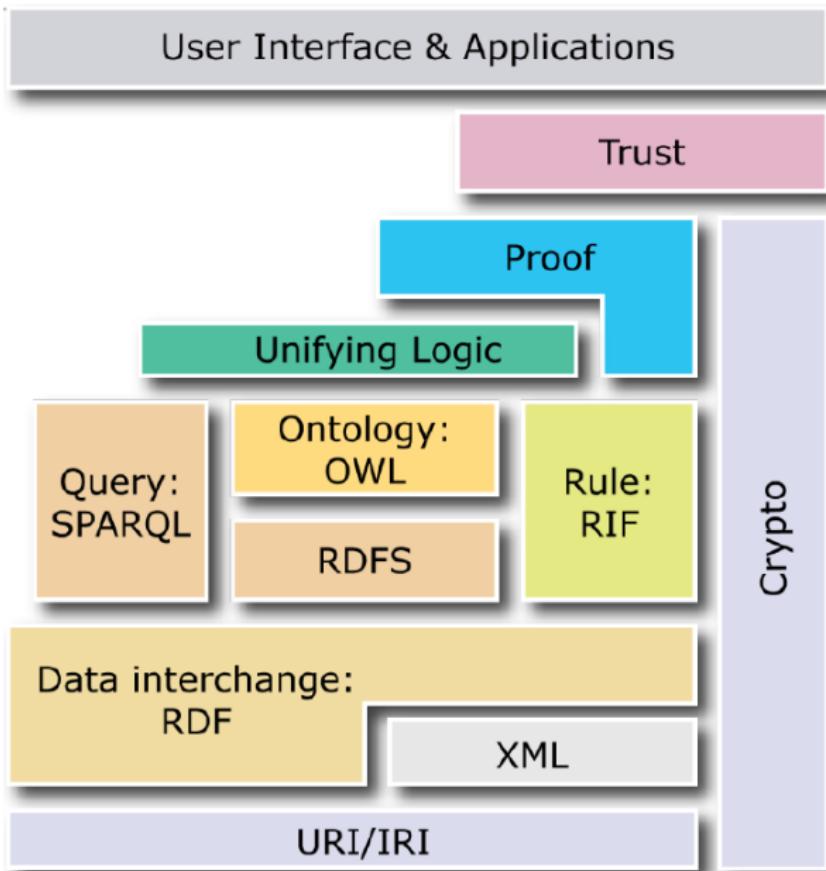
## 3 Example: Data Integration with Knowledge Graphs

## 4 Notes and Further Reading

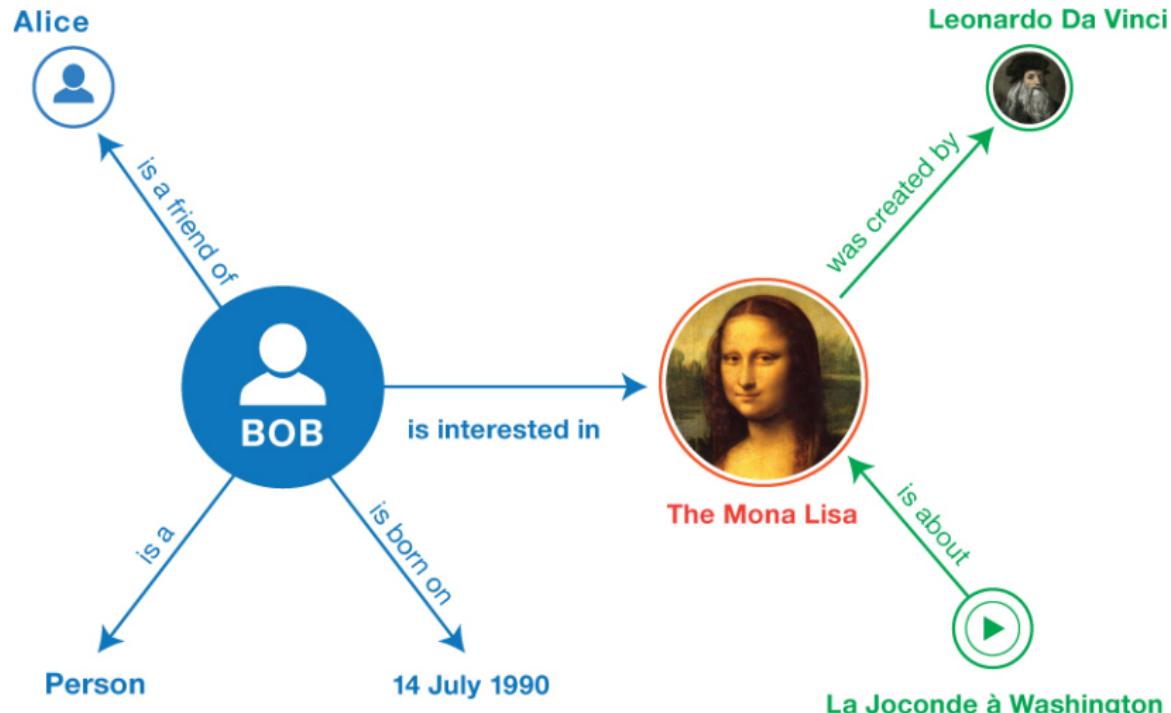
A dense network of neurons with glowing synapses.

The Basis: RDF

# The W3C “Layer Cake”



# Knowledge as Graphs



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

→ Worksheet #7: "Your first Knowledge Graph" & "Graph Updates"

# Triples

## Representation of Knowledge Graphs

In a system, we represent graphs in form of **triples**:

<subject> <predicate> <object>

(The *predicate* is sometimes called *property*.)

# Triples

## Representation of Knowledge Graphs

In a system, we represent graphs in form of **triples**:

<subject> <predicate> <object>

(The *predicate* is sometimes called *property*.)

## Examples

<*Bob*> <*is a*> <*person*>.

<*Bob*> <*is a friend of*> <*Alice*>.

<*Bob*> <*is born on*> <*the 14th of July 1990*>.

<*Bob*> <*is interested in*> <*the Mona Lisa*>.

<*the Mona Lisa*> <*was created by*> <*Leonardo da Vinci*>.

→ Worksheet #7: "Triples"

# Graphs vs. Triples

<subject> <predicate> <object>



→ Worksheet #7: "More Triples"

# RDF Triples

## The Resource Description Framework (RDF)

W3C (World Wide Web Consortium) standard ("recommendation")

- first public draft 1997
- RDF 1.0 in 1999; revised in 2004
- RDF 1.1 in 2014 (current version)

Family of standards: RDF, RDFS, RDFa, Turtle, N3, SPARQL, ...

# RDF Triples

## Format of triples

In RDF,

- Subject and predicate must be URIs (IRIs)
- Object can be IRI or [literal](#)

# RDF Triples

## Format of triples

In RDF,

- Subject and predicate must be URIs (IRIs)
- Object can be IRI or **literal**

## Examples

```
<http://www.wikidata.org/entity/Q12418>
<http://purl.org/dc/terms/title>
"Mona Lisa" .
```

```
<http://www.wikidata.org/entity/Q12418>
<http://purl.org/dc/terms/creator>
<http://dbpedia.org/resource/Leonardo_da_Vinci> .
```

→ Worksheet #7: "Wikidata" & "Using URIs"

# RDF Literals

"Mona Lisa"

In this triple:

```
<http://www.wikidata.org/entity/Q12418>
  <http://purl.org/dc/terms/title> "Mona Lisa" .
```

"Mona Lisa" is a **string literal**

# RDF Literals

"Mona Lisa"

In this triple:

```
<http://www.wikidata.org/entity/Q12418>
  <http://purl.org/dc/terms/title> "Mona Lisa" .
```

"Mona Lisa" is a **string literal**

## Things to know about literals

- Literals have a **datatype**, e.g., **string** or **int**
- Strings can have a **language tag**, e.g.,  
  "**Leonardo da Vinci**"@en  
  "**Léonard de Vinci**"@fr
- Strings are often used to provide human-readable **labels**
- For strings **only**, datatype can be omitted:  
  "**Mona Lisa**" is equivalent to "**Mona Lisa**"^^xsd:string
- Again, literals can **only** appear in the **object** position of a triple

All the details about datatypes:

<https://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/#section-Datatypes>

# Namespaces

## Shortening URIs

Instead of always writing full URIs (IRIs), we can split them into a **prefix** and **suffix**, e.g.:

```
<http://dbpedia.org/resource/Leonardo_da_Vinci>
```

- We define a prefix `dbpedia`:

```
PREFIX dbpedia: <http://dbpedia.org/resource/>
```

- and now we can simply write:

```
dbpedia:Leonardo_da_Vinci
```

- Note: angle brackets `<>` only for full IRIs

→ reduces dataset sizes, easier to read

# Namespaces

## Shortening URIs

Instead of always writing full URIs (IRIs), we can split them into a **prefix** and **suffix**, e.g.:

<[http://dbpedia.org/resource/Leonardo\\_da\\_Vinci](http://dbpedia.org/resource/Leonardo_da_Vinci)>

- We define a prefix **dbpedia**:

PREFIX dbpedia: <<http://dbpedia.org/resource/>>

- and now we can simply write:

dbpedia:Leonardo\_da\_Vinci

- Note: angle brackets <> only for full IRIs

→ reduces dataset sizes, easier to read

## Conventions

Commonly used URLs use the same namespace prefix

- E.g., FOAF (friend-of-a-friend):

PREFIX foaf: <<http://xmlns.com/foaf/0.1/>>

- Lookup a prefix at <https://prefix.cc/>

→ Worksheet #7: "More URIs" & "Namespaces"

# Serialization

## Formats

There is no single format .rdf (like .xml), commonly used are:

RDF/XML for data exchange (somewhat deprecated)

RDFa for embedding RDF into web pages

N-Triples (N3) for streaming RDF data and bulk dataset up-/download

Turtle for human-readable files

JSON-LD for web applications

plus some variations/extensions.

# Serialization

## Formats

There is no single format .rdf (like .xml), commonly used are:

RDF/XML for data exchange (somewhat deprecated)

RDFa for embedding RDF into web pages

N-Triples (N3) for streaming RDF data and bulk dataset up-/download

Turtle for human-readable files

JSON-LD for web applications

plus some variations/extensions.

## N-Triples

So far, we've mostly used the N-Triples format:

```
<http://www.wikidata.org/entity/Q12418> ←  
<http://purl.org/dc/terms/title> "Mona Lisa"
```

each line in a file is one triple, full IRIs only (no namespace prefixes)  
and ended by a period '.'

# Turtle

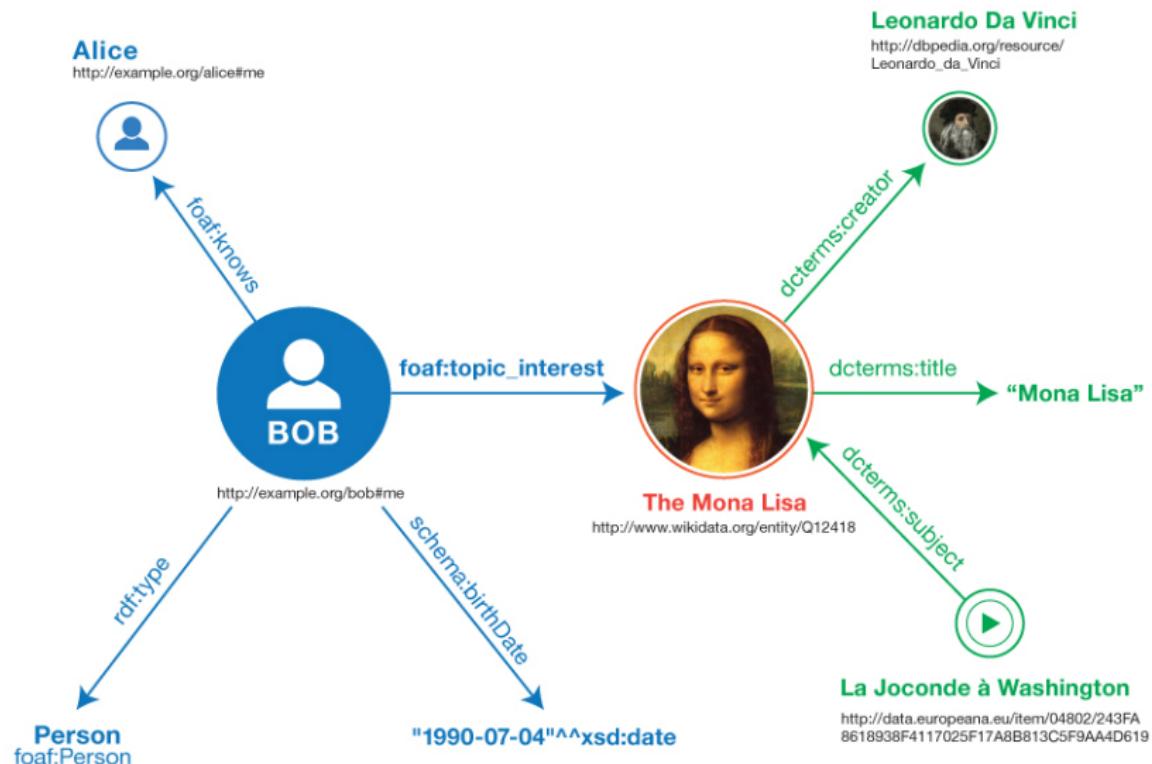
```
BASE  <http://example.org/>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema>
PREFIX schema: <http://schema.org/>
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX wd: <http://www.wikidata.org/entity/>

<bob#me>
    a foaf:Person ;
    foaf:knows <alice#me> ;
    schema:birthDate "1990-07-04"^^xsd:date ;
    foaf:topic_interest wd:Q12418 .

wd:Q12418
    dcterms:title "Mona\u201cLisa" ;
    dcterms:creator <http://dbpedia.org/resource/Leonardo_da_Vinci> .

<http://data.europeana.eu/item/04802/243FA8618938F4117025F17A8B813C5F9A
    dcterms:subject wd:Q12418 .
```

# Graph corresponding to the Turtle example



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

# RDF in programming practice

---

- ▶ For example, using Python+RDFLib:
  - a “Graph” object is created
  - the RDF file is parsed and results stored in the Graph
  - the Graph offers methods to retrieve (or add):
    - triples
    - (property,object) pairs for a specific subject
    - (subject,property) pairs for specific object
    - etc.
  - the rest is conventional programming...
- ▶ Similar tools exist in Java, PHP, etc.

# Python example using RDFLib

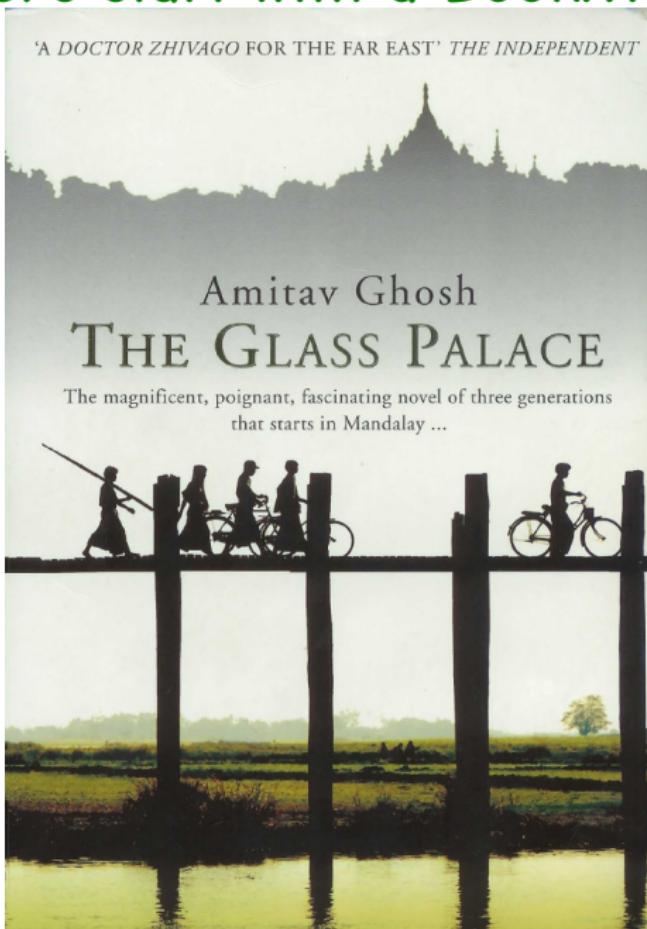
---

```
# create a graph from a file
graph = rdflib.Graph()
graph.parse("filename.rdf", format="rdfformat")
# take subject with a known URI
subject = rdflib.URIRef("URI_of_Subject")
# process all properties and objects for this subject
for (s,p,o) in graph.triples((subject,None,None)) :
    do_something(p,o)
```

# Outline

- 1 Introduction
- 2 The Resource Description Framework (RDF)
- 3 Example: Data Integration with Knowledge Graphs
- 4 Notes and Further Reading

# Example: Let's start with a Book...



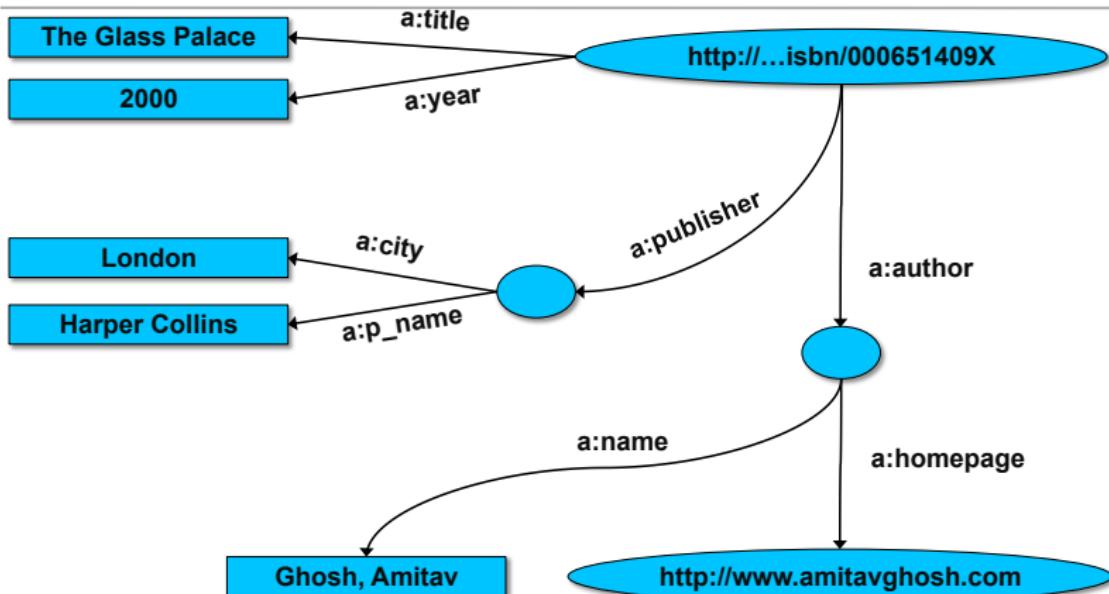
# A simplified bookstore data (dataset “A”)

ISBN	Author	Title	Publisher	Year
0006511409X	id_xyz	The Glass Palace	id_qpr	2000

ID	Name	Homepage
id_xyz	Ghosh, Amitav	<a href="http://www.amitavghosh.com">http://www.amitavghosh.com</a>

ID	Publisher's name	City
id_qpr	Harper Collins	London

# 1<sup>st</sup>: export your data as a set of relations

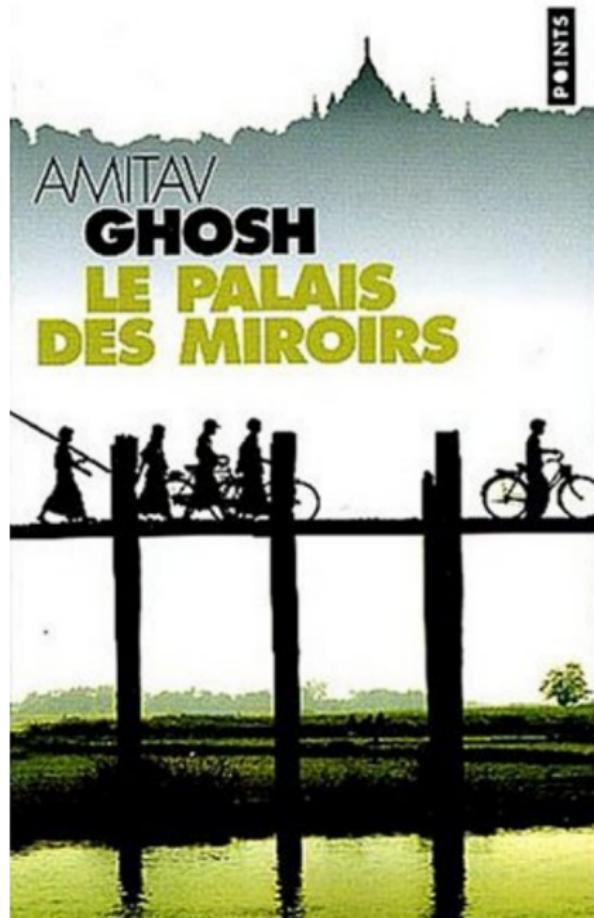


# Some notes on the exporting the data

---

- ▶ Relations form a graph
  - the nodes refer to the “real” data or contain some literal
  - how the graph is represented in machine is immaterial for now

Now the same book in French...

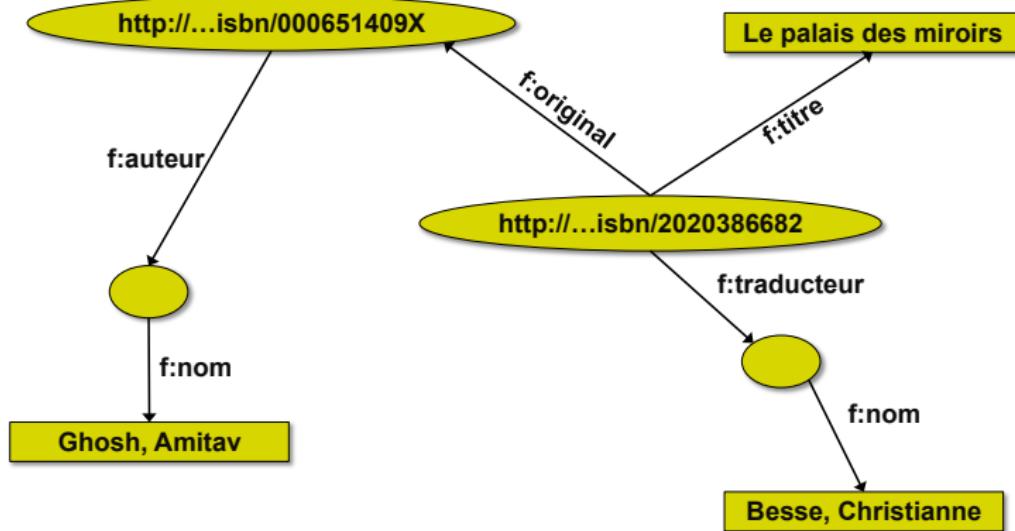


# Another bookstore data (dataset “F”)

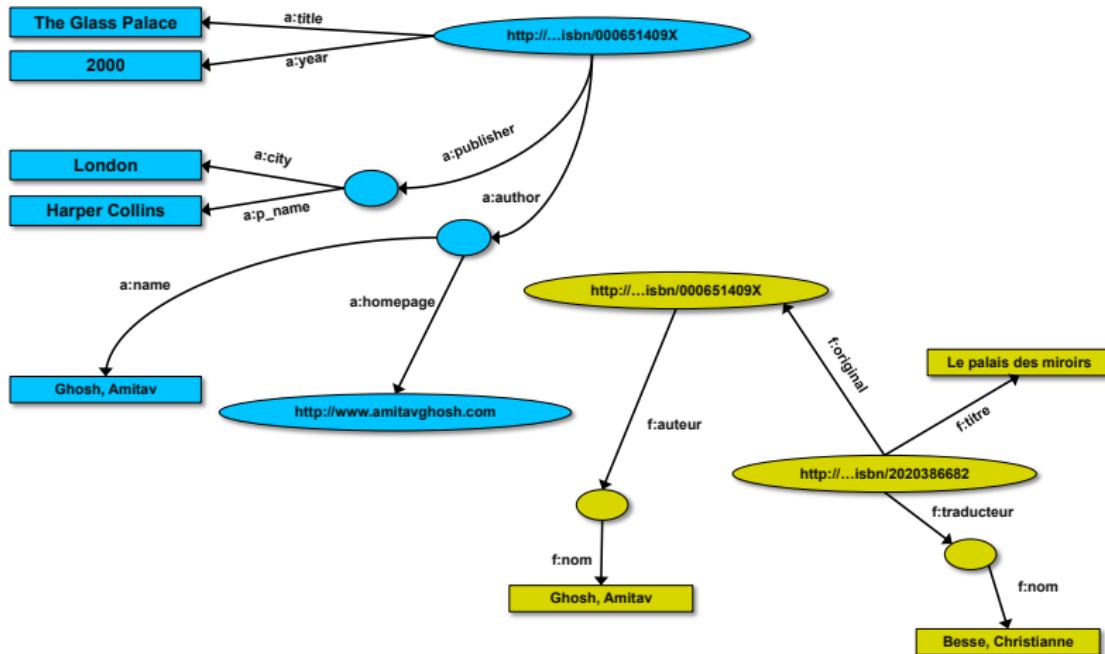
---

A	B	C	D
1	<b>ID</b>	<b>Titre</b>	<b>Traducteur</b>
2	ISBN 2020286682	Le Palais des Miroirs	\$A12\$
3			ISBN 0-00-6511409-X
4			
5			
6	<b>ID</b>	<b>Auteur</b>	
7	ISBN 0-00-6511409-X	\$A11\$	
8			
9			
10	<b>Nom</b>		
11	Ghosh, Amitav		
12	Besse, Christianne		

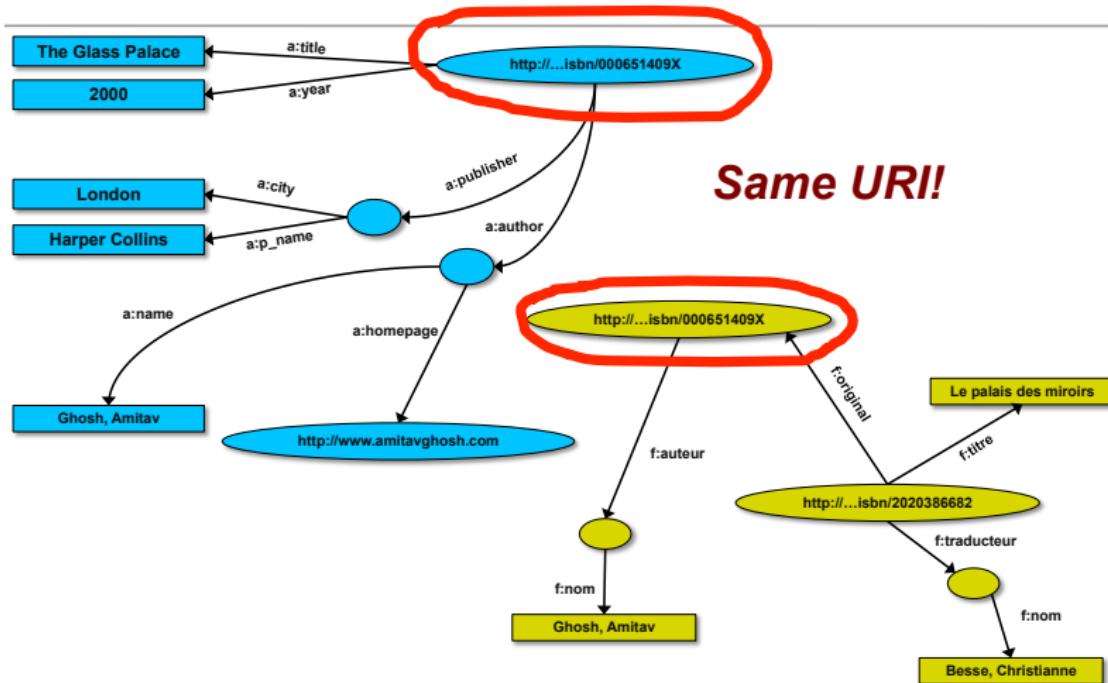
## 2<sup>nd</sup>: export your second set of data



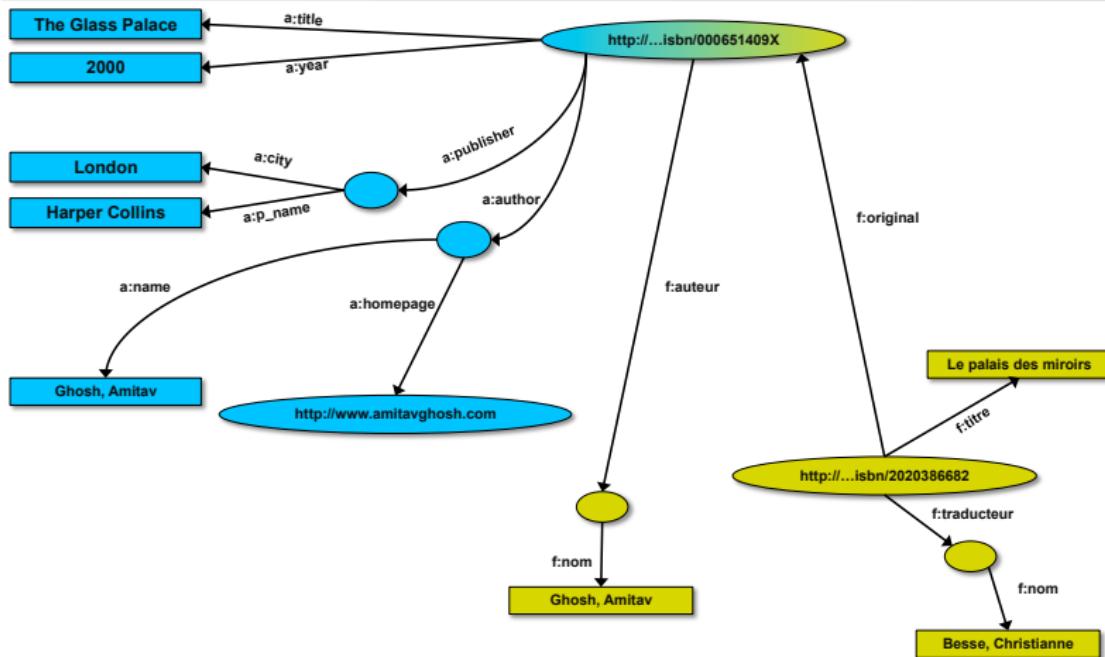
### 3<sup>rd</sup>: start merging your data



### 3<sup>rd</sup>: start merging your data (cont)



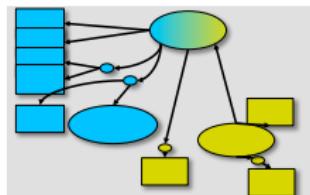
### 3<sup>rd</sup>: start merging your data



# Start making queries...

---

- ▶ User of data “F” can now ask queries like:
  - “give me the title of the original”
    - well, ... « donne-moi le titre de l’original »
- ▶ This information is not in the dataset “F”...
- ▶ ...but can be retrieved by merging with dataset “A”!

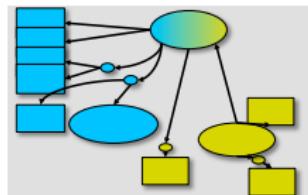
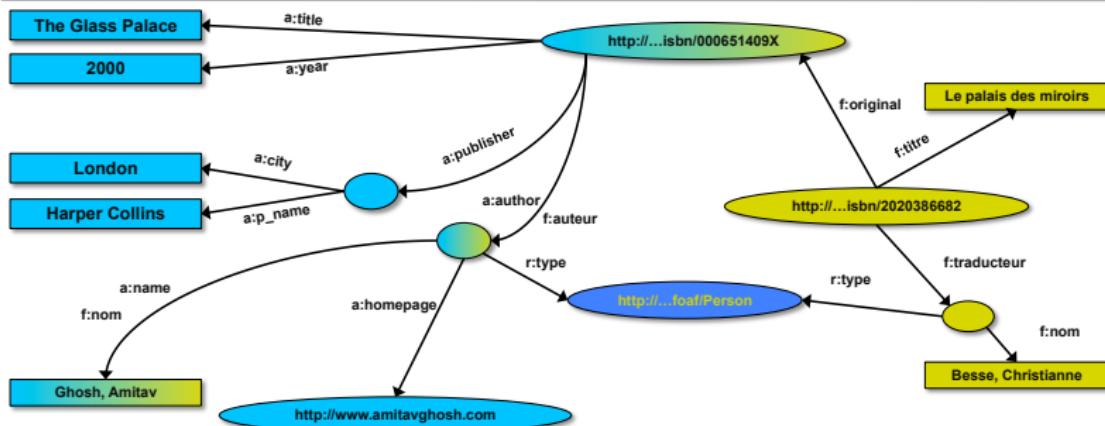


# However, more can be achieved...

---

- ▶ We “feel” that a:author and f:auteur should be the same
- ▶ But an automatic merge does not know that!
- ▶ Let us add some extra information to the merged data:
  - a:author same as f:auteur
  - both identify a “Person”
  - a term that a community may have already defined:
    - a “Person” is uniquely identified by his/her name and, say, homepage
    - it can be used as a “category” for certain type of resources

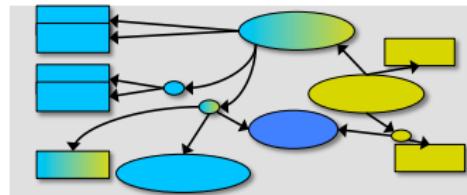
# 3<sup>rd</sup> revisited: use the extra knowledge



# Start making richer queries!

---

- ▶ User of dataset “F” can now query:
  - “donnes-moi la page d'accueil de l'auteur de l'original”
    - well... “give me the home page of the original's ‘auteur’”
- ▶ The information is not in datasets “F” or “A”...
- ▶ ...but was made available by:
  - merging datasets “A” and datasets “F”
  - adding three simple extra statements as an extra “glue”

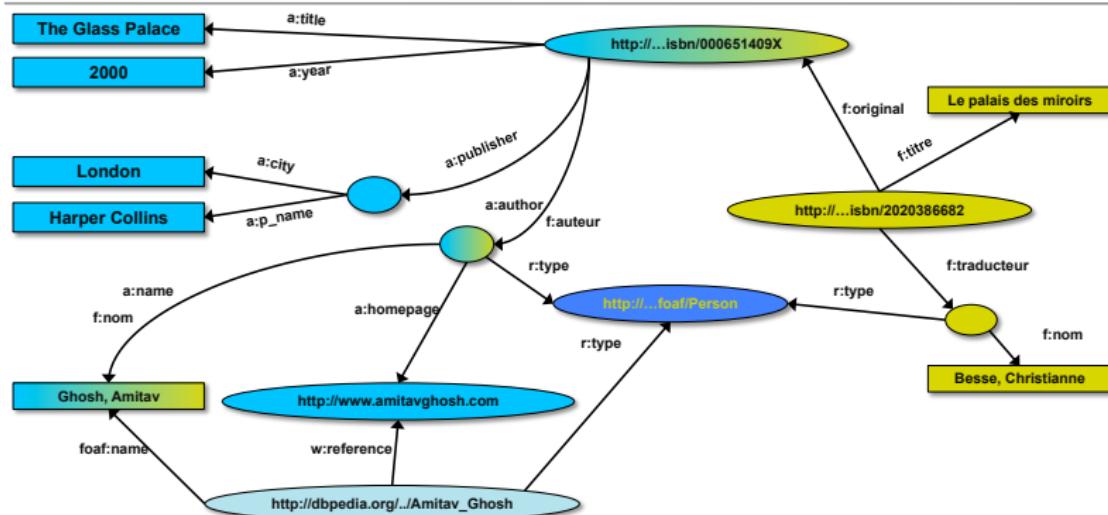


# Combine with different datasets

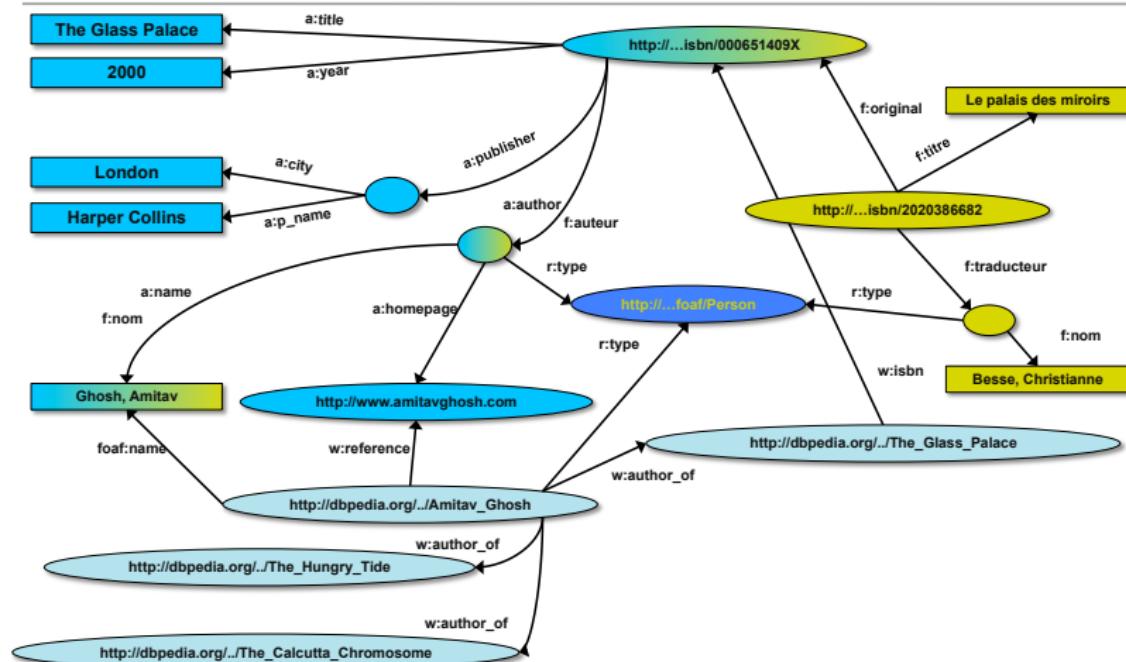
---

- ▶ Using, e.g., the “Person”, the dataset can be combined with other sources
- ▶ For example, data in Wikipedia can be extracted using dedicated tools
  - e.g., the “[dbpedia](#)” project can extract the “infobox” information from Wikipedia already...

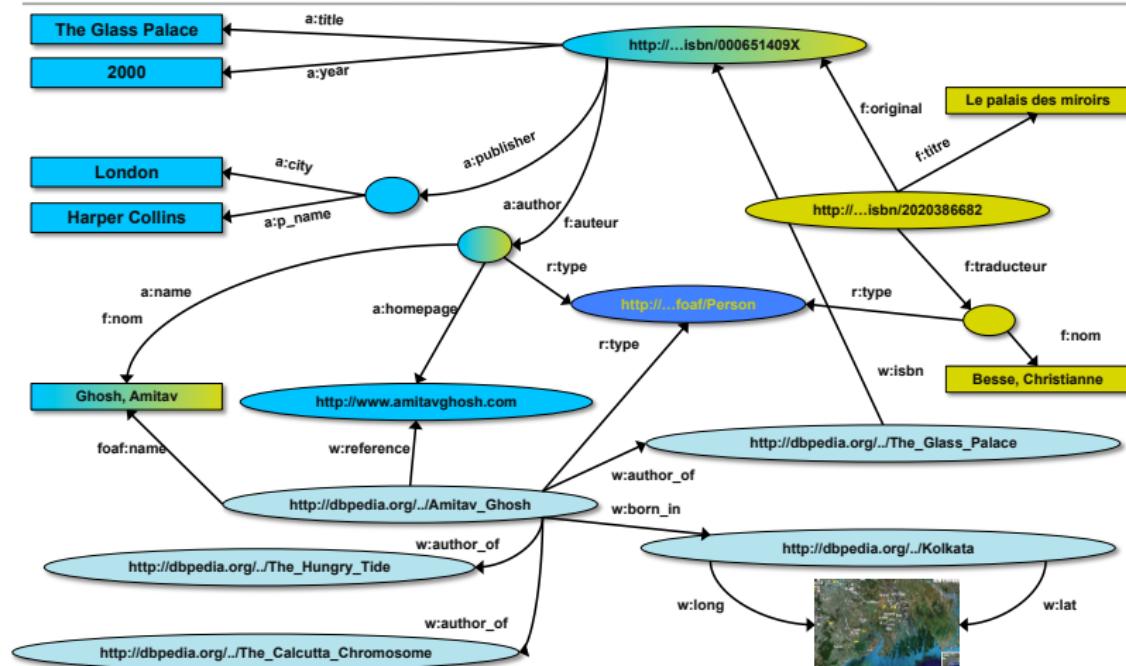
# Merge with Wikipedia data



# Merge with Wikipedia data



# Merge with Wikipedia data



# Is that surprising?

---

- ▶ It may look like it but, in fact, it should not be...
- ▶ What happened via automatic means is done every day by Web users!
- ▶ The difference: a bit of extra rigour so that machines could do this, too

# Outline

- 1 Introduction
- 2 The Resource Description Framework (RDF)
- 3 Example: Data Integration with Knowledge Graphs
- 4 Notes and Further Reading

# Reading Material

## Required

- [Wor14, Sections 1–3] (RDF Primer)

## Supplemental

- [Yu14, Chapters 1, 2] (Introduction, RDF)

# References

- [Her] Ivan Herman.  
Tutorial on Semantic Web Technologies.  
[http://www.w3.org/People/Ivan/CorePresentations/RDFTutorial/.](http://www.w3.org/People/Ivan/CorePresentations/RDFTutorial/)
- [Wor14] World Wide Web Consortium (W3C).  
RDF 1.1 Primer.  
<http://www.w3.org/TR/rdf11-primer/>, 24 June 2014.
- [Yu14] Liyang Yu.  
*A Developer's Guide to the Semantic Web.*  
Springer-Verlag Berlin Heidelberg, 2nd edition, 2014.  
<https://concordiauniversity.on.worldcat.org/oclc/897466408>.