# DIGITAL HUMANITIES SEMANTIC NOTEBOOKS

jupyter

```python
 1 occupation_list_nobles_status = ['No real information',
 2                                  'Soldiers',
 3                                  'Officials',
 4                                  'Self employed',
 5                                  'Craftsmen',
 6                                  'Merchants',
 7                                  'Retired',
 8                                  'Annuitants'
 9                                  ]
10 occupation_count_list = [occupation_nobles.count('Adel'),
11                          occupation_nobles.count('Militär'),
12                          occupation_nobles.count('Beamte'),
13                          occupation_nobles.count('Selbständig'),
14                          occupation_nobles.count('Handwerk'),
15                          occupation_nobles.count('Handel'),
16                          occupation_nobles.count('Rentner'),
17                          occupation_nobles.count('Rentier')
18                          ]
19 labels = occupation_list_nobles_status
20 sizes = occupation_count_list
21 explode = (0, 0, 0.1, 0, 0, 0, 0, 0)  # only "explode" the 2nd slice:'d
22 fig1, ax1 = plt.subplots()
23 ax1.pie(sizes, explode=explode, labels=labels, shadow=True, labeldistan
24 ax1.axis('equal')  # Equal aspect ratio ensures that pie is drawn as a
25 plt.title("German nobles' occupation in Paris in 1854")
26 plt.show()
```
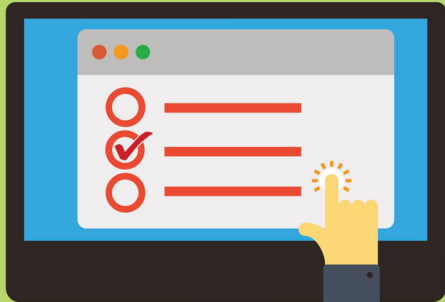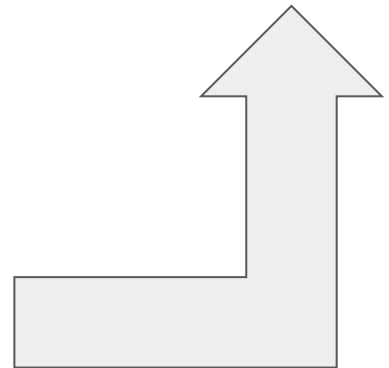
German nobles' occupation in Paris in 1854

No real information

Annuitants
Merchants
Craftsmen

Self employed

Officials

Soldiers

DHIP IHA

gkembellec@dhi-paris.fr

- Semantic publishing?
- Executable notebook?
                    (*aka* Jupyter Notebook)

https://pingo.coactum.de/337082

Please
fill in the form

# The challenges of science mediation

**John Unsworth**
**DH researcher**
**Dean of Virginia Libraries**

Scholarly
Primitives
(2001)
→
What do we need?

Semantic publishing
(2009)
→

How to do that ?

**David Shotton**
**Cambridge researcher**
**Semantic Publishing**
**Director of Open Citations**

Scholars' needs and ways to make linked reproducible
science

# What do we need as scholars ?

Scholarly Primitives, according to Unsworth (2000) :
- **Discovering** (serendipity)
- **Annotating** (Read / Write as a single process)
- **Comparing** (Theories, datasets, bibliographies)
- **Refering** (quote)
- **Sampling** (data, corpora…)
- **Illustrating** (examples, figures)
- **Representing**
        (disseminate in the community)

# 1st : semantic Web / publishing In science & culture

some meaningful concepts

# Semantic Web ?

Linking and describing resources
on the Web with :

- Authorities databases
  (ISNI, GnD, Wikidata, Dbpedia…)
- Vocabularies
  (schema.org, FoaF, Dublin Core...)
- A theoretical framework as a model
  (RDF)
- Technical implementations (TEI,
  HTML5 with RDFa, JSON-LD,



Sir Tim Berners Lee
Computer Scientist
Father of the Web (1991)
Father of the Semantic Web (2001)

# Semantic publishing as a solution ?

Semantic publishing best practices, according to David Shotton (2009) :
- Use established **standards** wherever possible.
- Publish **raw datasets** to the Web.
- Release **article metadata**, particularly reference lists, in **machine-readable form**.

# "Resource Description Framework"
## acronym RDF

RDF grammar of a triple :
( Subject , Predicate , Object )

Subject        : The resource
(what we are talking about) ;
Predicate    : the kind of resource property  ;
Object        : value of the property :
can be a number, a text, a web URL

# Document contents with RDF triples

If I write :

"<mark>Sociology of Religion has been written by Max Weber</mark>",

it could be seen as an RDF triple this way:

Subject          : ***"Sociology of Religion"***
Predicate        : ***was written by***
Object           : "***Max Weber***"

Using the description authority "Dublin Core" we will have "creator" for predicate :

**("Sociology of Religion"**, dc:creator , "Max Weber")

**http://purl.org/dc/elements/1.1/creator**          **same**

Description de ressources sur le Web

# Document contents with RDF triples

As there is an authority for books (ISBN) , we can use it to be sure we are talking about the good book :

Subject : ***isbn:***0-8070-4205-6
Predicate : http://purl.org/dc/elements/1.1/creator
Object : ***Max Weber***

Using the description authority "isbn" we will have 0-8070-4205-6 for subject

(***isbn:***0-8070-4205-6 ,
"http://purl.org/dc/elements/1.1/creator",
Max Weber)

# Document contents with RDF triples

There are also several authorities for people (ISNI, PND from GND, VIAF…) , we can use it to be sure we are talking about the good guy :

Subject            : *isbn:*0-8070-4205-6
Predicate          : http://purl.org/dc/elements/1.1/creator
Object             : http://d-nb.info/gnd/118629743

Using the description authority *Personennamendatei* or **PND** we will have "creator" for Object

(*isbn:*0-8070-4205-6,
http://purl.org/dc/elements/1.1/creator,
http://d-nb.info/gnd/118629743)

# Document contents with RDF triples : schema.org

lePack   Cnam   Pédagogie   Perso   CfP   Recherche   OpenDataSoft   stats   GestionDeTexte   Antidote   siteWeb   SEO   »   Autre

**Schema.org**        Documentation    Schemas    About

## Book
*A Schema.org Type*

Thing > CreativeWork > Book

[more...]

A book.

| Property | Expected Type | Description |
|---|---|---|
| **Properties from Book** | | |
| abridged | Boolean | Indicates whether the book is an abridged edition. |
| bookEdition | Text | The edition of the book. |
| bookFormat | BookFormatType | The format of the book. |
| illustrator | Person | The illustrator of the book. |
| isbn | Text | The ISBN of the book. |
| numberOfPages | Integer | The number of pages in the book. |
| **Properties from CreativeWork** | | |
| about | Thing | The subject matter of the content. Inverse property: subjectOf |
| abstract | Text | An abstract is a short description that summarizes a CreativeWork. |
| accessMode | Text | The human sensory perceptual system or cognitive faculty through which a person may process or perceive information. Expected values include: auditory, tactile, textual, visual, colorDependent, chartOnVisual, chemOnVisual, diagramOnVisual, mathOnVisual, musicOnVisual, textOnVisual. |

13

# Document contents with RDF triples

We can here define what the book is about using the book schema from schema.org/ :

Subject        : *"Sociology of Religion"*
Predicate    : https://schema.org/Book/About
Object        : "sociology", "religion"

# Semantic publishing?



Legend
Scientific paper
Formal semantics

No readable metadata
**Classical publishing**

Readable metadata
What has been called **semantic publishing**

Readable linked metadata
**Genuine semantic publishing**

Kuhn, T. and Dumontier, M.
'*Genuine Semantic Publishing*'
Data Science, vol. 1, no. 1-2,
pp. 139-154, 2017

# Small example : *mirodata* in a HTML page



Who's the "I"?

I do agree with *Max Weber's Sociology of Religion*

# Small example : *mirodata* in a HTML page



Subject      Predicate      Object

I do agree with *Max Weber's Sociology of Religion*

Who's the "I"?

Small example : *microdata* in a HTML page

```html
<!-- someone agrees -->
<p itemscope itemtype="http://schema.org/AgreeAction">
    I
    <!-- me, by the way, identified with my orcid id -->
    <span itemscope itemprop="agent"
          itemtype="http://schema.org/Person"
          itemid="https://orcid.org/0000-0003-3036-6989">
        <meta itemprop="name" content="Gérald Kembellec">
    </span>
do agree with
    <!-- with a book identified by its isbn id -->
<cite itemscope
       itemtype="http://schema.org/Book">
    <meta itemprop="isbn" content="0316769487">
    <span itemscope itemprop="author"
          itemtype="http://schema.org/Person"
          itemid="http://d-nb.info/gnd/118629743" >
        <span itemprop="name">Max Weber</span>'s
    </span>
    <span itemprop="name">
        Sociology of Religion
    </span>
</cite>
</p>
```

# Small example : *microdata* in a HTML page test with https://validator.schema.org/ 1/3

```html
1  <p itemscope itemtype="http://schema.org/AgreeAction">
2      I
3      <!-- me, by the way, identified with my orcid id -->
4      <span itemscope itemprop="agent"
5              itemtype="http://schema.org/Person"
6              itemid="https://orcid.org/0000-0003-3036-6989">
7              <meta itemprop="name" content="Gérald Kembellec">
8      </span>
9      do agree with
10     <!-- with a book identified by its isbn id -->
11     <cite itemscope
12             itemtype="http://schema.org/Book">
13             <meta itemprop="isbn" content="0316769487">
14             <span itemscope itemprop="author"
15                     itemtype="http://schema.org/Person"
16                     itemid="http://d-nb.info/gnd/118629743" >
17                     <span itemprop="name">Max Weber</span>'s
18             </span>
19             <span itemprop="name">
20                     Sociology of Religion
21             </span>
22     </cite>
23 </p>
```
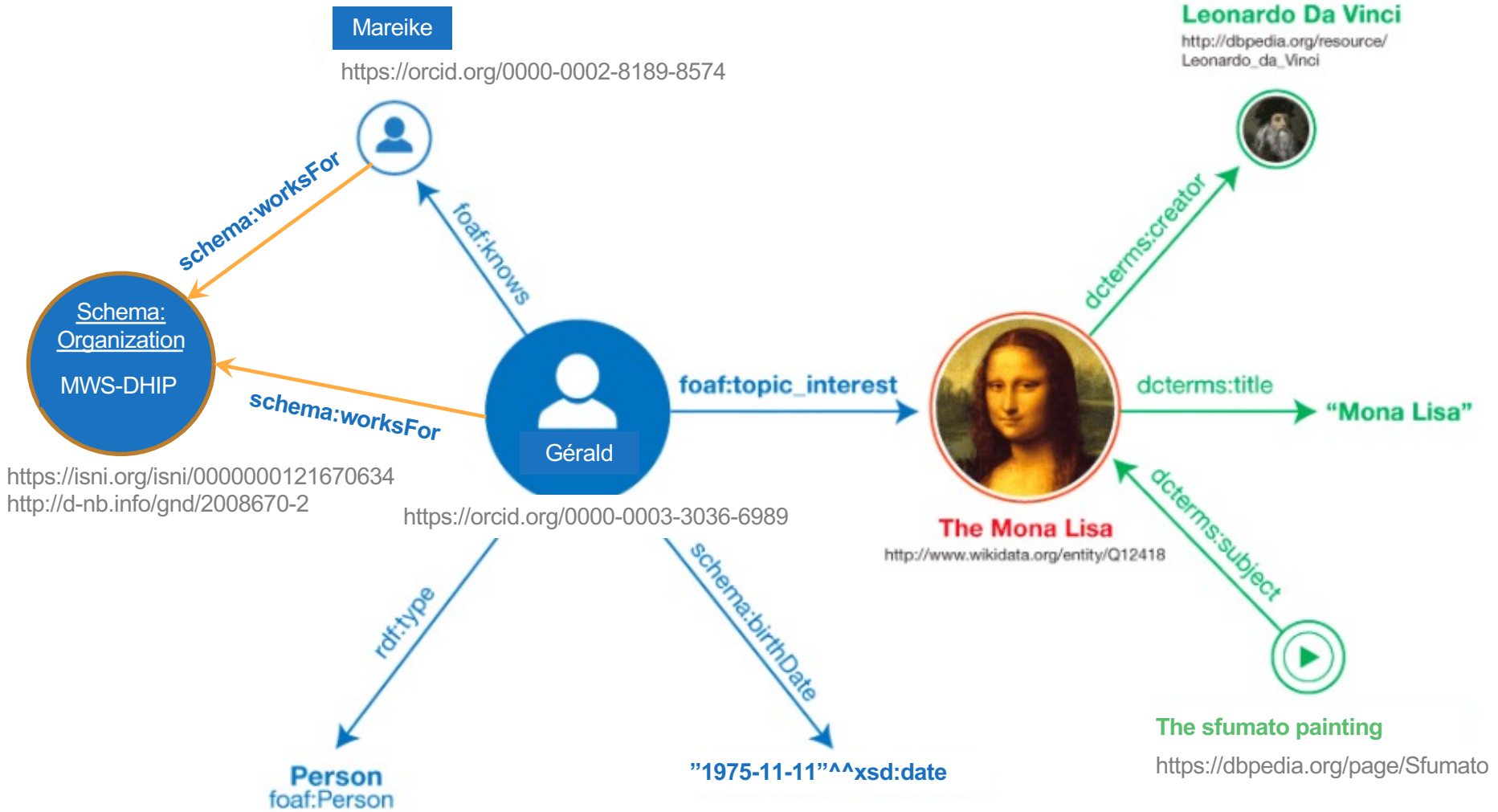
Détectés    0 ERREUR    0 AVERTISSEMENT    2 ÉLÉMENTS

AgreeAction    0 ERREUR    0 AVERTISSEMENT    1 ÉLÉMENT

Book    0 ERREUR    0 AVERTISSEMENT    1 ÉLÉMENT

# Small example : *microdata* in a HTML page
test with https://validator.schema.org/ 2/3

# Small example : *microdata* in a HTML page
## test with https://validator.schema.org/ 3/3

# Openlink Softwares used in the demonstration

openlink/**virtuoso-opensource**

Virtuoso is a high-performance and scalable Multi-Model RDBMS, Data Integration Middleware, Linked Data Deployment, and HTTP Application Server Platform

SparQL Enpoint : Virtuoso
https://virtuoso.openlinksw.com/



Browser's plugin
for semantic discovering tool :
Open Link structured data sniffer
https://addons.mozilla.org/fr/firefox/addon/openlink-structured-data-sniff/

# 2nd concept :
# The Jupyter Notebook

## What is a Jupyter notebook ?



This is a notebook



This is a laptop computer notebook



This is a Jupyter notebook

The first "notebook" by Galileo : manuscript from 17th century : "Observations of Jupiter  and four of its moons".

A "Jupyter notebook" is composed by :
- dataset pieces
      - raw / quite unstructured (txt, word…)
      - semi-structured (csv, json…)
      - structured (database)

- sections including
      - text (md, html…)
      - computed results (python, R…)
      - visual results (maps, graphs, charts…)

```
 1  #print(data.loc[(data['Beruf_Kategorie']=='noble')])
 2  number_chevaliers=data.value_counts(data['legion_d_honneur']=="Chevalier")
 3  number_officiers=data.value_counts(data['legion_d_honneur']=="Officier")
 4  number_grand_officiers=data.value_counts(data['legion_d_honneur']=="Grand Officier")
 5  number_grand_croix=data.value_counts(data['legion_d_honneur']=="Grand-Croix")
 6  print("In the dataset of noble German in Paris in 1854, dealing with 'la Légion d'honneur':")
 7  print("-",number_chevaliers.values[1],"had the 'chevalier' grade")
 8  print("-",number_officiers.values[1],"had the 'officier' grade")
 9  print("-",number_grand_officiers.values[1],"had 'grand officier' grade")
10  print("- and none had 'Commander' or 'Grand-Croix' grade")
```

```
In the dataset of noble German in Paris in 1854, dealing with 'la Légion d'honneur':
- 13 had the 'chevalier' grade
- 2 had the 'officier' grade
- 1 had 'grand officier' grade
- and none had 'Commander' or 'Grand-Croix' grade
```

A "Jupyter notebook" can be seen as a scientific production :
- with its own formalism
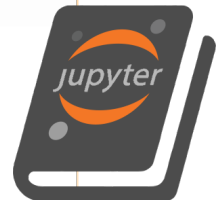- including

text

citations

footnotes, bibliographies

facts and figures

...

```
[ ]  1  #print(data.loc[(data['Beruf_Kategorie']==job),:])
     2  number_chevaliers=data.value_counts(data['legion_d_honneur']=="Chevalier")
     3  number_officiers=data.value_counts(data['legion_d_honneur']=="Officier")
     4  number_grand_officiers=data.value_counts(data['legion_d_honneur']=="Grand officier")
     5  number_grand_croix=data.value_counts(data['legion_d_honneur']=="Grand-Croix")
     6  print("In the dataset of noble German in Paris in 1854, dealing with 'la Légion d'honneur':")
     7  print("-",number_chevaliers.values[1],"had the 'chevalier' grade")
     8  print("-",number_officiers.values[1],"had the 'officier' grade")
     9  print("-",number_grand_officiers.values[1],"had 'grand officier' grade")
    10  print("- and none had 'Commander' or 'Grand-Croix' grade")
```

```
In the dataset of noble German in Paris in 1854, dealing with 'la Légion d'honneur':
- 13 had the 'chevalier' grade
- 2 had the 'officier' grade
- 1 had 'grand officier' grade
- and none had 'Commander' or 'Grand-Croix' grade
```

# What does it look like ?

```
Abbe von Jager, Civil, no_rank
Graf Miglied des Institus von Pradel Von, O., Civil, Chevalier
There were 19 nobles with 'Légion d'honneur' on the population of 116
```

```python
decorations_list_ldh = ['Chevalier','Officier','Grand Officier']
decorations_count_ldh = [decorations.count('Chevalier'),decorations.count('Officier'),decorations.count('Grand Officier')]
plt.bar(decorations_list_ldh,decorations_count_ldh)
plt.title("Number of german nobles in Paris with the 'légion d'honneur' in 1854")
plt.xlabel("Rank of 'Légion d'honneur'")
plt.ylabel("Number of units")
plt.show()
```

Number of german nobles in Paris with the 'légion d'honneur' in 1854



```python
percent_of_nobles_with_lgh = nb_lgh / noble_number * 100
print(str(round(percent_of_nobles_with_lgh)) + "% of noble Germans in Paris in 1854 received one of the 'Légion d'honneur' distinction")
```

```
16% of noble Germans in Paris in 1854 received one of the 'Légion d'honneur' distinction
```

```python
# Not the best code ever
soldier_status=data.value_counts(data['occupation_group']=="Militär")
noble_status=data.value_counts(data['occupation_group']=="Adel")
functionary_status=data.value_counts(data['occupation_group']=="Beamte")
self_employed_status=data.value_counts(data['occupation_group']=="Selbständig")
craftsmen_status=data.value_counts(data['occupation_group']=="Handwerk")
merchants_status=data.value_counts(data['occupation_group']=="Handel")
retired_status=data.value_counts(data['occupation_group']=="Rentner")
annuitant_status=data.value_counts(data['occupation_group']=="Rentier")
print("Most of nobles's activities are not clearly defined : "+ str(noble_status[True])+" are described with the 'Adel' term, which is quite fuzzy.")
print("The 1st category of occupation by ranking is 'functionaries' with a value of "+str(functionary_status[True])+" persons.")
# not significant :
```

<u>Ok, what do we need to do that ?</u>
- A notebook infrastructure
  - https://jupyter.org/try
  - https://jupyter-cloud.gwdg.de/ (DE)
  - https://colab.research.google.com/

- A research dataset (txt, json, geojson, csv…)
- A programming language (Python, R, Shell…)
- A presentation language (HTML or markdown)
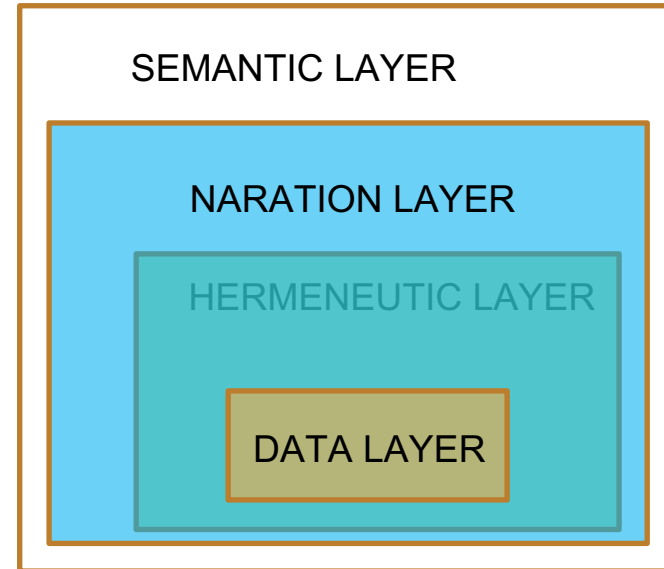- A metadata exposing language (COinS, RDFa, microdata…)
- Lot of work…

# Does it exists for real in the history field ?

**Journal of Digital History**

C²DH

DE G

SEMANTIC LAYER

NARATION LAYER

HERMENEUTIC LAYER

DATA LAYER

https://www.c2dh.uni.lu/news/digital-history-journal-transmedia-storytelling-hermeneutics-data

**Let's take a tour on Adressbuch1854's project**

… and go for demonstration…

… then we will discuss "semantic publishing" and "notebooks"

gkembellec@dhi-paris.fr

31