

Connecting a French Dictionary from the Beginning of the 20th Century to Wikidata

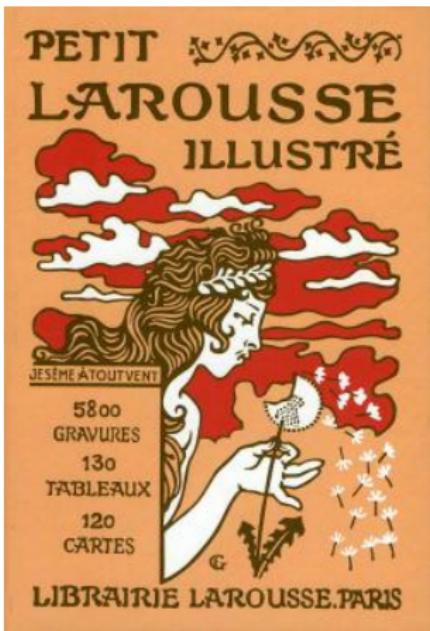
Pierre Nugues

Pierre.Nugues@cs.lth.se

LREC 2022, Marseilles, June 21-23, 2022

Petit Larousse illustré (1905)

- The *Petit Larousse illustré* is a one-volume dictionary of French, first published in 1905.
- Definitions of words, things, and people complemented with encyclopedic developments
- Numerous illustrations
- Very popular and far-reaching cultural influence.
- Updated yearly

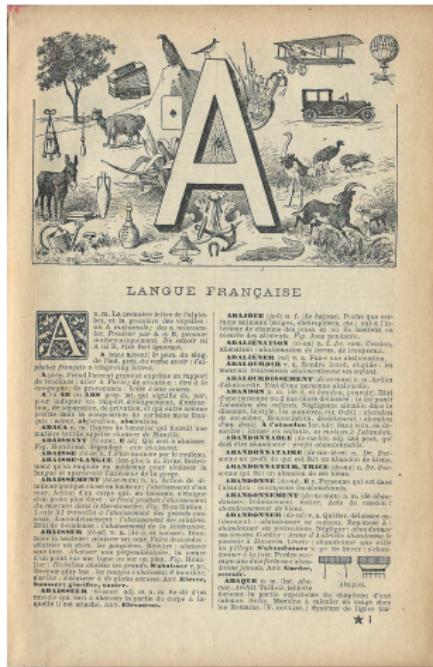


Source of image: Wikimedia

Dictionary Division

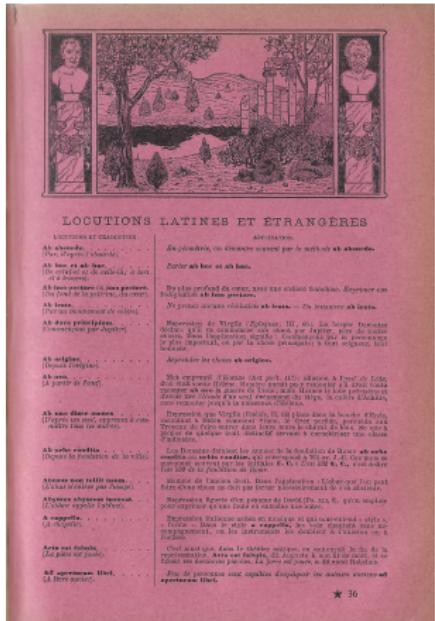
The dictionary consists of three main parts:

- ① The first one, *langue française*, is a dictionary of French with some encyclopedic content. It is restricted to the common nouns, verbs, adjectives, adverbs, and grammatical words (1066 pages);



4th edition from 1925

② The second part, *locutions*, contains quotes, mostly from classical Latin authors. It is much smaller than the two other parts (32 pages); and



4th edition from 1925

Dictionary Division

- ③ The third part, *histoire et géographie*, contains short encyclopedic descriptions of people, countries, locations, intellectual or art works (660 pages).



4th edition from 1925

Core Cultural Knowledge

- The *Larousse* content is highly didactic and corresponds to a core cultural knowledge in France at the beginning of the 20th century
- It is a snapshot of the world's view and a repository of the cultural values at that time.
- The printed content of such a dictionary is then highly informative
- It is nonetheless limited and static by nature.
- Elementary information is sometimes missing or incomplete:
 - Places of birth and death of persons;
 - Gender and occupations;
 - Exact geographical coordinates of locations; etc.
- Makes the (automatic) analysis sometimes difficult.

Entity Linking

- In this work, we connected manually all the 20,245 entries from the *history and geography* part to wikipedia information (wikidata items).
- Enables us to:
 - Extract more data on the entities;
 - Verify information;
 - Apply some automatic analyses;
 - Compare historically-situated representations with currently available descriptions;
- The annotated dataset is available in the JSON format from https://github.com/pnugues/petit_larousse_1905/.

- As reference points and identifiers, we used the wikidata items.
- Wikidata assigns a unique identification number to each thing in wikipedia
- Pierre Larousse has identifier Q313709,
- Wikidata was designed to handle the multilingual versions of wikipedia.
- Contains structured information

Item Discussion Read View history

Pierre Larousse (Q313709)

French grammarian, lexicographer and encyclopaedist
Pierre Athanase Larousse

In more languages
Català

Language	Label	Description
English	Pierre Larousse	French grammarian, lexicographer and encyclopaedist
French	Pierre Larousse	grammaireien, lexicographe et encyclopédiste français
Italian	Pierre Larousse	linguista, pedagogista e editore francese
Swedish	Pierre Larousse	No description defined

All entered languages

Statements

instance of	human
	2 references

image



- As reference points and identifiers, we used the wikidata items.
- Wikidata assigns a unique identification number to each thing in wikipedia
- Pierre Larousse has identifier Q313709,
- Wikidata was designed to handle the multilingual versions of wikipedia.
- Contains structured information

Wikipedia (35 entries)

ar	بيير لارويس
arz	بيير لارويس
az	Pyer Laruss
bg	Пиер Ларус
br	Pierre Larousse
ca	Pierre Larousse
cs	Pierre Athanase Larousse
da	Pierre Larousse
de	Pierre Larousse
el	Πιέρ Λαρούς
en	Pierre Larousse
eo	Pierre Larousse
es	Pierre Larousse
eu	Pierre Larousse
fi	Pierre Larousse
fr	Pierre Larousse
hr	Pierre Larousse

- As reference points and identifiers, we used the wikidata items.
- Wikidata assigns a unique identification number to each thing in wikipedia
- Pierre Larousse has identifier Q313709,
- Wikidata was designed to handle the multilingual versions of wikipedia.
- Contains structured information

date of birth	23 October 1817 <i>Gregorian</i>
	► 12 references
place of birth	Toucy
	► 2 references
date of death	3 January 1875 <i>Gregorian</i>
	► 11 references
place of death	Paris
	► 2 references
cause of death	cerebrovascular disease
	► 1 reference
place of burial	Montparnasse Cemetery
	► 0 references

Method

- We manually connected each entry of the *Larousse* to one or more wikidata items.
- We started from the OCRed and corrected version from the Nénufar website.
- Nénufar provides an XML-TEI structure of the entries, but we only used the raw text

<http://nenufar.huma-num.fr/>



Nénufar

Le Petit Larousse illustré de 1906 à 1948 - Suivre l'ajout des données.

Edition Toutes Choisir Article Votre recherche Plus de critères...

Langue A B C D E F G H I J K L M N O P Q R S T U V W X Y Z - Locutions - Noms propres A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

1687 articles

AALI-PACHA,

homme d'Etat turc, né à Constantinople. Il a attaché son nom à la politique de réformes du Tanzimat (1815-1871).

Example of Link

Some entries are easy to link like the person names:

AALI-PACHA, homme d'Etat turc, né à Constantinople. Il a attaché son nom à la politique de réformes du Tanzimat (1815-1871).

'AALI-PASHA, Turkish statesman, born in Constantinople. He attached his name to the policy of reforms of the Tanzimat (1815-1871).'

corresponds to Q439237 as the occupation and the dates of birth and death match those in wikidata.

given name	Mehmed » 0 references
	Emin » 0 references
	Ali » 0 references
noble title	pasha » 0 references
date of birth	5 March 1815 Gregorian » 4 references
place of birth	Istanbul » 1 reference
date of death	7 September 1871 Gregorian » 3 references



Example of Link

Some other entries, like countries, can change in nature and form.

*GRANDE-BRETAGNE et
IRLANDE (Royaume-Uni
de)...*

*'GREAT BRITAIN and IRE-
LAND (United Kingdom
of)...*'

We chose the Q174193 identifier, which describes a state that existed from 1801 to 1927 when the Larousse was published.

United Kingdom of Great Britain and Ireland (Q174193)

historical state (1801–1922), name in use until 1927

United Kingdom | UK | GB | UKGBI | Great Britain and Ireland | the United Kingdom | Britain

[+ In more languages](#)

[Configure](#)

Language	Label	Description	Also known as
English	United Kingdom of Great Britain and Ireland	historical state (1801–1922), name in use until 1927	United Kingdom UK GB UKGBI Great Britain and Ireland the United Kingdom Britain
French	Royaume-Uni de Grande-Bretagne et d'Irlande	État historique (1801–1922), nom en usage jusqu'en 1927	
Italian	Regno Unito di Gran Bretagna e Irlanda	stato europeo esistito dal 1801 al 1922	Regno di Gran Bretagna Regno Unito di Gran Bret
Swedish	Förenade kungariket Storbritannien och Irland	statsbildning mellan nuvarande Storbritannien och Irland 1801–1922	Förenade kungariket Storbritannien och Irland Förenade konungariket S

All entered languages

Statements

instance of	sovereign state	edit
	start time	1 January 1801 Gregorian
	end time	12 April 1927 Gregorian
	follows	Great Britain

Tricky Cases

Works of art sometimes have changed name or attribution.

For example, this painting in the *Petit Larousse illustré*:

Automne (I'), tableau de Jordaeus (Bruxelles) ;

When ambiguous, we relied on the more detailed *Nouveau Larousse illustré* (1897–1904):

*Automne (REPRÉSENTATIONS DI-
VERSES DE L'). On connaît l'Automne,
tableau de Jordaeus (musée de Brux-
elles), allégorie des occupations et des
dons de l'automne. (Cette déesse,
drapée dans un manteau rouge, a les
mains pleines de raisins. Parmi les per-
sonnages qui l'entourent, on remarque
une nymphe nue, vue de dos, aux formes
opulentes) ;*



Jacob Jordaens - Allegory of Fertility - Google Art Project, Q22787456

Entries with Two or More Links

Some entries contain lists of people or works of art, as for example:

ABAD Ier [bad'], premier roi maure de Séville, et chef de la dynastie des Abadites ; il régna de 1023 à 1042. — Son fils Abad II régna de 1042 à 1069, et son petit-fils, Abad III, de 1069 à 1095.

'ABBAD I, first Moorish king of Seville, and head of the Abbadid dynasty; he reigned from 1023 to 1042. — His son Abbad II reigned from 1042 to 1069, and his grandson, Abbad III, from 1069 to 1095.'

- Most of the time, such lists contain a dash character ‘–’ to separate the entities.
- Using this character, we extracted all the lists and we connected all the names mentioned in them.
- For the previous example, we linked the three entities to:
["Q305795", "Q30556", "Q299578"] .

Dataset Format and Size

- We stored the annotation as a list of dictionaries in a JSON file.
- Each dictionary represents an entry and has two keys: the raw text of the entry and a list of wikidata identifiers, most often only one.
- For the AALI-PASHA entry, this corresponds to:

```
"texte": "AALI-PACHA, homme d'Etat turc, né à  
Constantinople. Il a attaché son nom à la  
politique de réformes du Tanzimat  
(1815-1871).",  
"qid": ["Q439237"]
```

- At the end, the complete annotation of the 20,245 entries resulted in 22,357 links, where 18,905 entries have one link and 1340 have two or more.
- The annotated dataset is available from GitHub:
https://github.com/pnugues/petit_larousse_1905/.

The Semantic Categories

- Using the identifiers, we extracted the semantic category of the entries from wikidata.
- We used the SPARQL query language and triples like this one:

```
wd:Q439237 wdt:P31 ?type .
```

where wd:Q439237 is Aali-Pasha's identifier, wdt:P31 is a property meaning *instance of*, and ?type is the type we want to extract.

- The SPARQL server returns the Q5 identifier denoting a human.
- When an entry contained more than one name, we considered only the first identifier in the list.

The Ten most Frequent Categories

Q-Number	Frequency	Description
Q5	7653	human
Q484170	3022	commune of France
Q1549591	849	big city
Q515	789	city
Q7725634	572	literary work
Q4022	447	river
Q3305213	308	painting
Q747074	255	comune of Italy
Q1637706	231	city with a population of more than 1,000,000
Q23442	221	island

The two first types are human with nearly 38% of the entries and French commune with nearly 15%.

The Human Beings

- Human being is the most frequent category and the corresponding wikidata items often provide the dates and places of birth and death of the person, sometimes the occupations, spouses, etc.
- Note that less known entities are not always completely documented.
- We extracted the dates of birth and death from wikidata.
- We extracted the dates with the SPARQL triples:

`wd:Q439237 wdt:P569 ?db .`

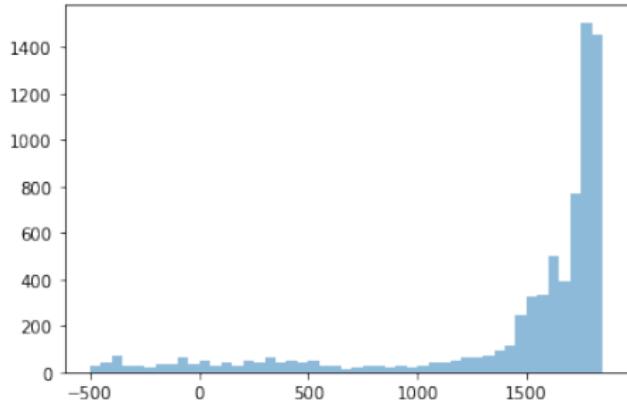
`wd:Q439237 wdt:P570 ?dd .`

where `wdt:P569` is the property for the date of birth and `wdt:P570` for the date of death.

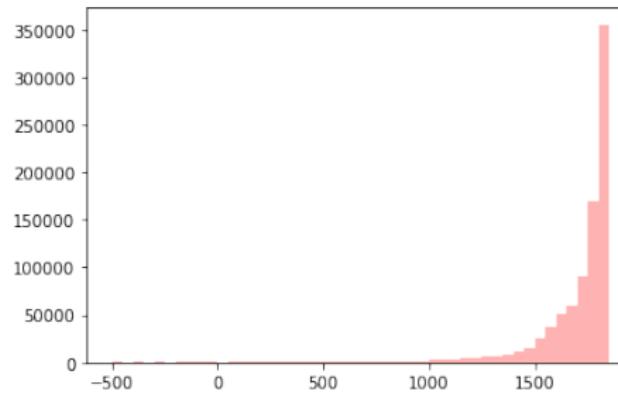
- The wikidata queries for the human beings returned 7430 pairs.

Distribution Across the Years: Birth Dates

- Excluding the dates before -500 and after 1851, we computed the distribution of the birth dates of the humans.



Petit Larousse illustré

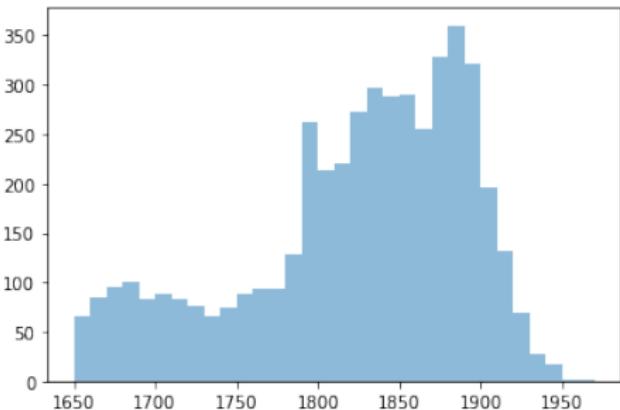


Wikidata

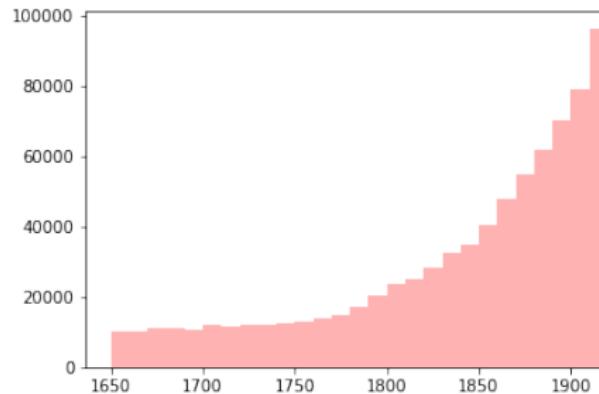
- The two largest bins are between 1750 and 1800 and 1800 and 1850. These two most recent bins are twice as large as the one just before them between 1700 and 1750.

Distribution Across the Years: Death Dates

- The death dates with a distribution starting at year 1650.



Petit Larousse illustré

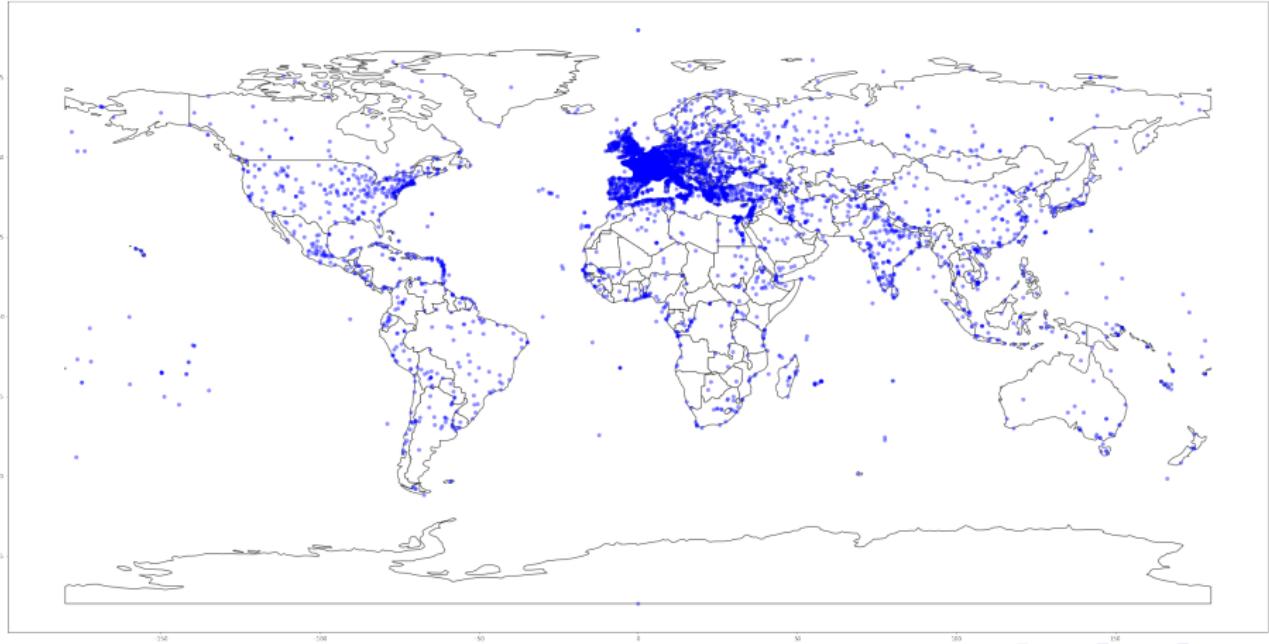


Wikidata

- Confirms the accent on contemporary or near-contemporary people. At the date of publication (1905), 334 people were still living
- Year 1800 is a turning point, reflecting a specific attention of the *Larousse* on the new era started by the French Revolution of 1789.

Geographical Entities

- The geographical entities in the *Larousse* convey a world view.
- We extracted their coordinates with the P625 property:
?e wdt:P625 ?geo .



Conclusion

- We annotated all the entries of the *history and geography* part of the *Petit Larousse illustré* from 1905 with their wikidata identifiers.
- We have outlined how to use the identifiers to extract, complement, and process knowledge on the entities.
- This has enabled us to exhibit more precisely the scope of information and world view conveyed by this dictionary.
- We have released the complete annotated corpus
- We hope that this language resource will be useful to train entity linking applications or facilitate new projects in digital humanities.