

TEI 2024

From Catullus to Wikidata: Language Models, Metadata Schemas, and Ontologies in a Digital Edition in TEI

Carlos Javier Nusch ¹, Gabriel Alejandro Calarco ², Gimena del Rio Riande ², Leticia Cecilia Cagnina ^{2,3}, Marcelo Luis Errecalde ^{2,3}, Leandro Antonelli ^{4,5}

¹PREBI-SEDICI, UNLP, ²CONICET, ³LIDIC, UNSL,
⁴LIFIA, UNLP, ⁵CAETI, UAI, Argentina.

carlosnusch@prebi.unlp.edu.ar



Esta obra está bajo una Licencia Creative Commons
Atribución-NoComercial-CompartirIgual 4.0 Internacional



Aetatis Amoris Project: Markup and NLP for Latin Love Poetry

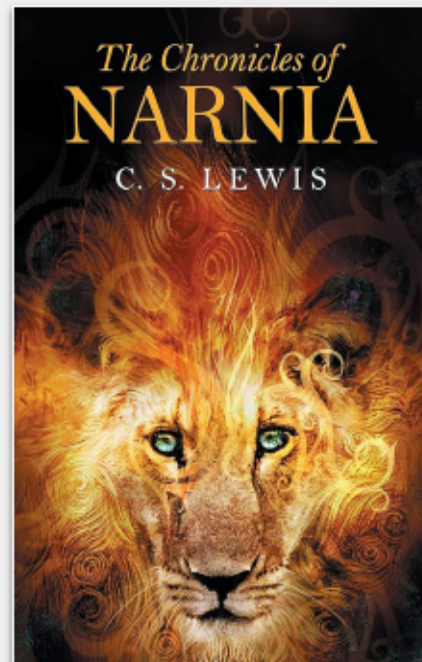
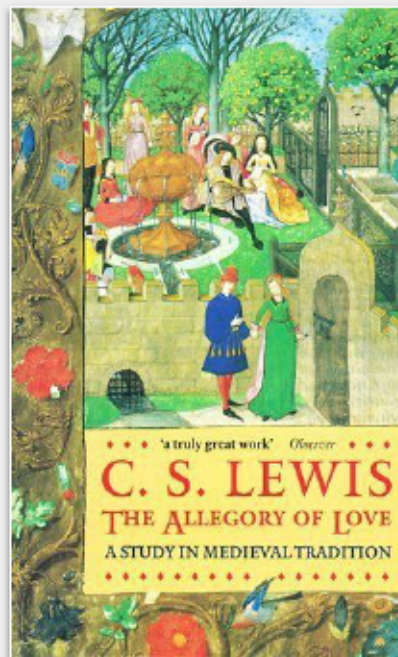
- This presentation describes the markup and **natural language processing** tasks undertaken in the *Aetatis Amoris* project, which explores love poetry and its textual patterns throughout literary history.
- In this early stage, the project focuses on the works of classical Latin poets such as *Catullus*, *Tibullus*, and *Propertius*.
- One approach involves using the **TEI-XML** standard to encode these texts.



Motivation: A Reading of C. S. Lewis



- “French poets, in the eleventh century, discovered or invented, or were the first to express, that romantic species of passion which English poets were still writing about in the nineteenth.” (Lewis, 1936)
- Lewis suggested that the shape of the love motifs and imaginary we had, at least at the beginnings of the XX Century, belong to the occitan poets in the XI century.



A few minutes editing tasks

The first tasks involved automatically generating **key document elements** like the **header** and **body**, along with counting and **automatically tagging verses** and **stanzas**.

```
1 from xml.etree.ElementTree import Element, SubElement, ElementTree
2 from pathlib import Path
3
4 # Rutas a las carpetas de entrada y salida
5 folder_paths = [
6     '/content/gdrive/My Drive/A Doctorado/datos_thesis/Catulo',
7     '/content/gdrive/My Drive/A Doctorado/datos_thesis/Tibulo',
8     '/content/gdrive/My Drive/A Doctorado/datos_thesis/Propercio',
9 ]
10 output_folder = Path('/content/gdrive/My Drive/A Doctorado/datos_thesis/Output_TEI')
11
12 # Detalles de la edición original para el encabezado TEI
13 editions = {
14     'Catulo': ("Catullus; edited by Elmer Truesdell Merrill", "1893", "http://archive.org/details/catulluseditedby00catuuoft"),
15     'Tibulo': ("Tibulli aliorumque carminum libri tres. Scriptorum classicorum bibliotheca Oxoniensis", "1915", "https://archive.org/details/tibullialiorumqu00tibuoft"),
16     'Propercio': ("Sex. Propertii Elegiae", "1898", "https://archive.org/details/sexpropertiieleg00prop")
17 }
```

A few minutes editing tasks

The first tasks involved automatically generating **key document elements** like the **header** and **body**, along with counting and **automatically tagging verses and stanzas**.

```
19 def create_header(poem_title, author, edition, year, url):
20     header = Element('teiHeader')
21     file_desc = SubElement(header, 'fileDesc')
22     title_stmt = SubElement(file_desc, 'titleStmt')
23     title = SubElement(title_stmt, 'title')
24     title.text = poem_title
25     author_el = SubElement(title_stmt, 'author')
26     author_el.text = author
27
28     publication_stmt = SubElement(file_desc, 'publicationStmt')
29     publisher = SubElement(publication_stmt, 'publisher')
30     publisher.text = "Digital Edition"
31
32     source_desc = SubElement(file_desc, 'sourceDesc')
33     bibl = SubElement(source_desc, 'bibl')
34     bibl.text = f"{edition}, {year}. Available at: {url}"
35
36     return header
37
```

A few minutes editing tasks

The first tasks involved automatically generating **key document elements** like the **header** and **body**, along with counting and **automatically tagging verses** and **stanzas**.

```
1 from xml.etree import ElementTree
2 from pathlib import Path
3
4 # Función para contar versos en el texto y actualizar el archivo TEI
5 def count_verses_and_update_tei(tei_file_path):
6     print(f"Procesando archivo: {tei_file_path}")
7     # Cargar el archivo TEI existente
8     tree = ElementTree.parse(tei_file_path)
9     root = tree.getroot()
10    body = root.find('.//body')
11
```

```
# Contar versos y estrofas
verses = text.split('\n') # Suponiendo que cada verso está separado por un salto de línea
stanza_count = 0
verse_count = 0
current_stanza = None

for verse in verses:
    if verse.strip(): # Asegurarse de que no es una línea vacía
        if current_stanza is None:
            current_stanza = ElementTree.SubElement(body, 'lg')
            stanza_count += 1
            print(f"Nueva estrofa creada. Total estrofas: {stanza_count}")
        verse_element = ElementTree.SubElement(current_stanza, 'l')
        verse_element.text = verse
        verse_count += 1
        print(f"Verso añadido: {verse}")
    else:
        current_stanza = None # Resetear la estrofa si hay una línea en blanco
        print("Línea en blanco encontrada, estrofa reseteada.")

# Eliminar el párrafo original después de procesar su contenido
body.remove(paragraph)

# Agregar información del recuento al body
count_info = ElementTree.SubElement(body, 'count')
count_info.text = f'Total versos: {verse_count}, Total stanzas: {stanza_count}'
print(f"Total de versos: {verse_count}, Total de estrofas: {stanza_count}")

# Guardar los cambios en el archivo
tree.write(tei_file_path, encoding='unicode', xml_declaration=True)
```

A few minutes editing tasks

The first tasks involved automatically generating **key document elements** like the **header** and **body**, along with counting and **automatically tagging verses and stanzas**.

```

Catulo_Carmen_068a: Bloc de notas
Archivo Edición Formato Ver Ayuda
Quod mihi fortuna casuque oppressus acerbo
conscriptum hoc lacrimis mittis epistolium,
naufragum ut eiectum spumantibus aequoris undis
sublevem et a mortis limine restituum,
quem neque sancta Venus molli requiescere somno
desertum in lecto caelibe perpetitur,
nec veterum dulci scriptorum carmine musae
oblectant, cum mens anxia pervigilat,
id gratum est mihi, me quoniam tibi dicis amicum
muneraque et Musarum hinc petis et Veneris.
sed tibi ne mea sint ignota incommoda, Manli,
neu me odisse putes hospitis officium,
accipe quis merser fortunae fluctibus ipse,
ne amplius a misero dona beata petas.
tempore quo primum vestis mihi tradita pura est,
iucundum cum aetas florida ver ageret,
multa satis lusi; non est dea nescia nostri
quae dulcem curis miscet amaritiam:
sed totum hoc studium luctu fraterna mihi mors
abstulit. o misero frater adempte mihi,
tu mea tu moriens fregisti commoda, frater,
tecum una tota est nostra sepulta donus,
  
```

Línea 1, columna 1 | 100% | Windows (CRLF) | UTF-8

```

<TEI>
  <<teiHeader>
    <<fileDesc>
      <<titleStmt>
        <title>Catulo_Carmen_068a</title>
        <author>Catulo</author>
      </titleStmt>
      <<publicationStmt>
        <publisher>Digital Edition</publisher>
      </publicationStmt>
      <<sourceDesc>
        <bibl>Catullus; edited by Elmer Truesdell Merrill, 1893. Available at: http://archive.org/details/catulluseditedby00catuuoft</bibl>
      </sourceDesc>
    </fileDesc>
  </teiHeader>
  <<text>
    <<cooq>
      <<lg type="TIPO_DE_VERSO">
        <<lg type="TIPO_DE_ESTROFA" n="1">
          <l n="1">Quod mihi fortuna casuque oppressus acerbo</l>
          <l n="2">conscriptum hoc lacrimis mittis epistolium,</l>
          <l n="3">naufragum ut eiectum spumantibus aequoris undis</l>
          <l n="4">sublevem et a mortis limine restituum,</l>
        </lg>
        <<lg type="TIPO_DE_ESTROFA" n="2">
          <l n="5">quem neque sancta <persName key="Venus">Venus</persName> molli requiescere somno</l>
          <l n="6">desertum in lecto caelibe perpetitur,</l>
          <l n="7">nec veterum dulci scriptorum carmine musae</l>
          <l n="8">oblectant, cum mens anxia pervigilat,</l>
        </lg>
        <<lg type="TIPO_DE_ESTROFA" n="3">
          <l n="9">id gratum est mihi, me quoniam tibi dicis amicum</l>
          <l n="10">muneraque et Musarum hinc petis et Veneris.</l>
          <l n="11">sed tibi ne mea sint ignota incommoda, <persName key="Manli">Manli</persName>,</l>
          <l n="12">neu me odisse putes hospitis officium,</l>
        </lg>
      </cooq>
    </text>
  </TEI>
  
```

2th Task: Automatic Extraction and Tagging of Entities

Subsequently, using the advanced **spaCy model** for Latin, **LatinCy**, the names of people and places were automatically extracted and **tagged** with the corresponding labels.

```
!pip install spacy
```

```
!pip install https://huggingface.co/latincy/la\_core\_web\_lg/resolve/main/la\_core\_web\_lg-any-py3-none-any.whl
```



LatinCy Community

✕ diyclassics  diyclassics

Synthetic trained spaCy pipelines for Latin NLP

Developed by [Patrick J. Burns](#), 2023.

2th Task: Automatic Extraction and Tagging of Entities

Subsequently, using the advanced **spaCy model** for Latin, **LatinCy**, the names of people and places were automatically extracted and **tagged** with the corresponding labels.

```
1 import spacy
2 from pathlib import Path
3
4 # Cargar el modelo de spaCy
5 nlp = spacy.load("la_core_web_lg")
6
7 # Rutas a las carpetas de entrada
8 folder_paths = [
9     '/content/gdrive/My Drive/A_Doctorado/datos_thesis/Catulo',
10    '/content/gdrive/My Drive/A_Doctorado/datos_thesis/Tibulo',
11    '/content/gdrive/My Drive/A_Doctorado/datos_thesis/Propercio',
12 ]
13
14 # Listas para almacenar entidades
15 person_entities = []
16 location_entities = []
17
```



LatinCy Community

✕ diyclassics  diyclassics

Synthetic trained spaCy pipelines for Latin NLP

Developed by [Patrick J. Burns](#), 2023.

Automatic Extraction and Tagging of Entities

```
18 # Procesar cada archivo en las carpetas
19 for folder_path in folder_paths:
20     folder = Path(folder_path)
21     print(f"Procesando archivos en {folder.name}:")
22     for file_path in folder.glob('*.txt'):
23         with file_path.open('r', encoding='utf-8') as file:
24             text = file.read()
25             doc = nlp(text)
26             print(f"Entidades en {file_path.name}:")
27             for ent in doc.ents:
28                 if ent.label_ == "PERSON":
29                     person_entities.append(ent.text)
30                 elif ent.label_ == "LOC":
31                     location_entities.append(ent.text)
32                 print(f"{ent.text} ({ent.label_})")
33             print("-" * 30)
34
35 # Guardar las entidades en archivos
36 with open('/content/gdrive/My Drive/A_Doctorado/datos_thesis/person_entities.txt', 'w', encoding='utf-8') as f:
37     for person in person_entities:
38         f.write(person + '\n')
39
40 with open('/content/gdrive/My Drive/A_Doctorado/datos_thesis/location_entities.txt', 'w', encoding='utf-8') as f:
41     for location in location_entities:
42         f.write(location + '\n')
```

Automatic Extraction and Tagging of Entities

Subsequently, using the advanced spaCy model for Latin, LatinCy, names of people and places were automatically extracted and tagged with the corresponding labels.

```
Procesando archivos en Catulo:  
Entidades en Catulo_Carmen_029.txt:  
Mamurram (PERSON)  
Britannia (LOC)  
Cinaede (PERSON)  
Romule (PERSON)  
Adoneus (PERSON)  
Romule (PERSON)  
Tagus (LOC)  
Galliae (LOC)  
Britanniae (LOC)  
-----  
Entidades en Catulo_Carmen_014b.txt:  
-----  
Entidades en Catulo_Carmen_012.txt:  
Marrucine (PERSON)  
Asini (PERSON)  
Pollioni (PERSON)  
Saetaba (PERSON)  
Fabullus (PERSON)  
Veranius (PERSON)  
Veraniolum (PERSON)  
Fabullum (PERSON)
```

```
Procesando archivos en Catulo:  
Entidades en Catulo_Carmen_029.txt:  
Mamurram (PERSON)  
Britannia (LOC)  
Cinaede (PERSON)  
Romule (PERSON)  
Adoneus (PERSON)  
Romule (PERSON)  
Tagus (LOC)  
Galliae (LOC)  
Britanniae (LOC)
```

```
Entidades en Propercio_Liber_3_Poema_20.txt:  
Africa (LOC)  
Luna (LOC)  
Venus (PERSON)  
Amor (PERSON)
```

```
Entidades en Propercio_Liber_1_Poema_13.txt:  
Galle (LOC)  
Galle (LOC)  
Galle (LOC)  
Enipeo (PERSON)  
Taenarius (PERSON)  
Iove (PERSON)  
Ledae (PERSON)  
Ledae (PERSON)
```

Automatic Extraction and Tagging of Entities

Subsequently, using the advanced spaCy model for Latin, LatinCy, names of people and places were automatically extracted and tagged with the corresponding labels.

```
Procesando archivos en Catulo:  
Entidades en Catulo_Carmen_029.txt:
```

```
Mamurram (PERSON)  
Britannia (LOC)  
Cinaede (PERSON)  
Romule (PERSON)  
Adoneus (PERSON)  
Romule (PERSON)  
Tagus (LOC)  
Galliae (LOC)  
Britanniae (LOC)
```

```
-----  
Entidades en Catulo_Carmen_014b.txt:
```

```
-----  
Entidades en Catulo_Carmen_012.txt:
```

```
Marrucine (PERSON)  
Asini (PERSON)  
Pollioni (PERSON)  
Saetaba (PERSON)  
Fabullus (PERSON)  
Veranius (PERSON)  
Veraniolum (PERSON)  
Fabullum (PERSON)
```

```
Procesando archivos en Catulo:
```

```
Entidades en Catulo_Carmen_029.txt:
```

```
Mamurram (PERSON)  
Britannia (LOC)  
Cinaede (PERSON)  
Romule (PERSON)  
Adoneus (PERSON)  
Romule (PERSON)  
Tagus (LOC)  
Galliae (LOC)  
Britanniae (LOC)
```

```
Entidades en Propercio_Liber_3_Poema_20.txt:
```

```
Africa (LOC)  
Luna (LOC)  
Venus (PERSON)  
Amor (PERSON)
```

```
Entidades en Propercio_Liber_1_Poema_13.txt:
```

```
Galle (LOC)  
Galle (LOC)  
Galle (LOC)  
Enipeo (PERSON)  
Taenarius (PERSON)  
Iove (PERSON)  
Ledae (PERSON)  
Ledae (PERSON)
```

It was perfect until... The Metaphorical Use Paradox

What happens when gods are used **metaphorically** to represent places? For example, Uranus for heaven or Tethys and Oceanus for the sea or ocean? There is our first editorial decision to make...

Thetis (PERSON)
Nereine (PERSON)
Oceanus (PERSON)

Uraniae (LOC)



Automatic Extraction and Tagging of Entities

```
1 import spacy
2 from pathlib import Path
3 from lxml import etree
4
5 # Cargar el modelo de spaCy
6 nlp = spacy.load("la_core_web_lg")
7
8 # Función para crear el elemento TEI correcto
9 def create_tei_element(entity):
10     if entity.label_ == "PERSON":
11         return f'<persName key="{entity.text}">{entity.text}</persName>' #SACAR ATRIBUTO KEY
12     elif entity.label_ == "LOC":
13         return f'<placeName key="{entity.text}">{entity.text}</placeName>'
14     return entity.text
15
16 # Rutas a las carpetas de entrada
17 folder_paths = [
18     '"/content/gdrive/My Drive/A__Doctorado/datos_thesis/Output_TEI'
19 ]
20
```

Automatic Extraction and Tagging of Entities

```
21 # Procesar cada archivo TEI en las carpetas
22 for folder_path in folder_paths:
23     folder = Path(folder_path)
24     print(f"Procesando archivos en {folder.name}:")
25     for file_path in folder.glob('*.xml'): # Controlar la extensión correcta para los archivos TEI
26         tree = etree.parse(str(file_path))
27         root = tree.getroot()
28         # Procesar cada línea de poesía dentro de los elementos <l>
29         for l in root.xpath('//text//body//lg//l'):
30             original_text = l.text
31             doc = nlp(original_text)
32             new_text = original_text
33             for ent in reversed(doc.ents): # Reversed para no alterar los índices mientras reemplazamos
34                 start = ent.start_char
35                 end = ent.end_char
36                 tei_tag = create_tei_element(ent)
37                 new_text[:start] + tei_tag + new_text[end:]
38                 print(f"Reemplazando '{ent.text}' por '{tei_tag}' en línea {l.get('n')}")
39
40             # Actualizar el texto de la línea
41             l.text = new_text
```

Automatic Extraction and Tagging of Entities

Procesando archivos en Output_TEI:

```
Reemplazando 'Uraniae' por '<placeName key="Uraniae">Uraniae</placeName>' en línea 2
Reemplazando 'Manilo' por '<persName key="Manilo">Manilo</persName>' en línea 16
Reemplazando 'Vinia' por '<persName key="Vinia">Vinia</persName>' en línea 16
Reemplazando 'Phrygium' por '<placeName key="Phrygium">Phrygium</placeName>' en línea 18
Reemplazando 'Asia' por '<placeName key="Asia">Asia</placeName>' en línea 22
Reemplazando 'Thespiae' por '<placeName key="Thespiae">Thespiae</placeName>' en línea 27
Reemplazando 'Aonios' por '<persName key="Aonios">Aonios</persName>' en línea 28
Reemplazando 'Aganippe' por '<persName key="Aganippe">Aganippe</persName>' en línea 30
Reemplazando 'Venus' por '<persName key="Venus">Venus</persName>' en línea 61
Reemplazando 'Au-' por '<persName key="Au-">Au-</persName>' en línea 83
Reemplazando 'Oceano' por '<placeName key="Oceano">Oceano</placeName>' en línea 86
Reemplazando 'Talasio' por '<persName key="Talasio">Talasio</persName>' en línea 129
Reemplazando 'Tyrio' por '<placeName key="Tyrio">Tyrio</placeName>' en línea 167
Reemplazando 'Venus' por '<persName key="Venus">Venus</persName>' en línea 193
Reemplazando 'Venus' por '<persName key="Venus">Venus</persName>' en línea 197
Reemplazando 'Africi' por '<placeName key="Africi">Africi</placeName>' en línea 201
Reemplazando 'Torquatus' por '<persName key="Torquatus">Torquatus</persName>' en línea 211
Reemplazando 'Manlio' por '<persName key="Manlio">Manlio</persName>' en línea 217
Reemplazando 'Telemacho' por '<persName key="Telemacho">Telemacho</persName>' en línea 224
Archivo actualizado: Catulo_Catulo_Carmen_061_TEI.xml
Reemplazando 'Libystinis' por '<placeName key="Libystinis">Libystinis</placeName>' en línea 1
Reemplazando 'Scylla' por '<persName key="Scylla">Scylla</persName>' en línea 2
Archivo actualizado: Catulo_Catulo_Carmen_060_TEI.xml
Reemplazando 'Lesbia' por '<persName key="Lesbia">Lesbia</persName>' en línea 7
Reemplazando 'Catulle' por '<persName key="Catulle">Catulle</persName>' en línea 13
Archivo actualizado: Catulo_Catulo_Carmen_051_TEI.xml
```

What else can I do to enhance my edition?



What else can I do to enhance my edition?

- In addition, searches for standardized metadata were carried out using external **APIs**, and consulting databases such as Virtual International Authority File (VIAF), Pleiades project, and Wikidata.
- This allowed for the **retrieval** of **standardized identifiers**, **rich** and **curated information**, and **images** of historical places and characters.



VIAF



APIs

- API stands for **Application Programming Interface**, which is a **set of rules and protocols** that allows different software applications to **communicate** with each other.
- APIs provide a way for developers to access data or services from external systems without needing to know the internal workings of those systems.



VI
AF



What is VIAF?

- VIAF, or **Virtual International Authority File**, is a system that links **authority records of names** (people, organizations) from national libraries and other institutions worldwide.
- It helps standardize name **identifiers** across **databases** and **provides unique, persistent identifiers for individuals**, making it easier to connect various information sources.



VIAF

VIAF API

- We **searched** for the VIAF (Virtual International Authority File) **IDs** of a **list of person names**.
- We used the VIAF API to find the matching records, and then saved the found results (including the name, VIAF ID, and link) into one file and the names that were not found into another.



```
1 import requests
2
3 def get_viaf_details(name):
4     url = f"http://viaf.org/viaf/AutoSuggest?query={name}"
5     try:
6         response = requests.get(url)
7         data = response.json()
8         if 'result' in data and data['result']:
9             result = {
10                 "term": data['result'][0]['term'],
11                 "viaf_id": data['result'][0]['viafid'],
12                 "link": f"https://viaf.org/viaf/{data['result'][0]['viafid']}"
13             }
14             print(f"Nombre buscado: {name}, Nombre encontrado: {result['term']}, {result['viaf_id']}, {result['link']}")
15             return result
16         else:
17             print(f"Nombre buscado: {name}, Not Found")
18             return None
19     except requests.RequestException as e:
20         print(f"Error en la solicitud para {name}: {e}")
21         return None
22
23 # Cargar la lista de nombres de personas desde un archivo y eliminar duplicados
24 with open('/content/gdrive/My Drive/A_Doctorado/datos_thesis/person_entities.txt', 'r', encoding='utf-8') as f:
25     person_entities = set(line.strip() for line in f if line.strip())
26
```

VIAF API

The code worked perfectly, but the **volume of data retrieved was enormous!!!!**

```
Nombre buscado: Amabo, Nombre encontrado: Amabokoboko (Rugby team), 159851024, https://viaf.org/viaf/159851024
Nombre buscado: Thamyrae, Not Found
Nombre buscado: Tarquini, Nombre encontrado: Tarquini, 291167802705917770162, https://viaf.org/viaf/291167802705917770162
Nombre buscado: Licini, Nombre encontrado: Licinio Refice, 1883-1954, 14959385, https://viaf.org/viaf/14959385
Nombre buscado: Acmen, Nombre encontrado: Acmenide, 9514159477862627990005, https://viaf.org/viaf/9514159477862627990005
Nombre buscado: Claudius, Nombre encontrado: Claudius Galenus, ca. 129-ca. 210/216, 44299175, https://viaf.org/viaf/44299175
Nombre buscado: Rhenus, Nombre encontrado: Rhenus, 1575-1642, 49288294, https://viaf.org/viaf/49288294
Nombre buscado: Risi, Nombre encontrado: Risi, Dino, 1916-2008, 19761937, https://viaf.org/viaf/19761937
Nombre buscado: Crannonis, Not Found
Nombre buscado: Marathum, Not Found
Nombre buscado: Artaciis, Not Found
Nombre buscado: Threicia, Not Found
Nombre buscado: Romule, Nombre encontrado: Romuleus, Bartholomaeus, -1588, 38006675, https://viaf.org/viaf/38006675
Nombre buscado: Nereine, Not Found
Nombre buscado: Danai, Nombre encontrado: Danai Gurira, 1978-, 274583701, https://viaf.org/viaf/274583701
Nombre buscado: Tappone, Nombre encontrado: Tapponen, Sakari, 1538163706790929421526, https://viaf.org/viaf/1538163706790929421526
```

Pleiades Project

- Is an open, collaborative resource providing **detailed information** about **ancient places**, primarily from the Greco-Roman world.
- It links ancient and modern geography using **digital mapping tools** and **open data** like **geographic information systems (GIS)**, allowing users to visualize the locations on a map and see how ancient geography relates to today's world.



Data Download from Pleiades

First, the data was downloaded from the Pleiades Downloads page, where various datasets are available in formats like CSV, GeoJSON, and KML.

The chosen datasets were:

- **Names CSV:** Contains information on attested names of ancient places, their romanized forms, and links to Pleiades records.
- **Locations CSV:** Includes geographical coordinates (points or polygons) of registered places.

PLEIADES

[Log in](#)[Home](#)[Places](#)[Credits](#)[Participate](#)[Blog](#)[Documentation](#)[Downloads](#)

You are here: [Home](#) → [Data for download](#)

Data for download

Creators: [Sean Gillies](#)

Contributors: [Brian Turner](#), [Tom Elliott](#)

Copyright © The Contributors. Sharing and remixing permitted under terms of the Creative Commons Attribution 3.0 License (cc-by).

Last modified Jan 30, 2024 06:22 AM

Search

[Advanced Search...](#)

Upcoming Events

Processing Data from Pleiades

- We **filtered entities** that match the predefined list of historical names.
- Checked for matches in both the **attested** and **romanized** forms.
- The filtered data was merged with **geographical information (points and polygons)**. Finally, the results are saved in a new CSV file and displayed.

```
1 import pandas as pd
2
3 # Cargar los datos desde los archivos CSV
4 df_names = pd.read_csv('/content/gdrive/My Drive/A_Doctorado/datos_thesis/Pleiades_CSVs/names.csv', usecols=[
5 | 'place_id', 'attested_form', 'romanized_form_1', 'romanized_form_2', 'romanized_form_3', 'uri'
6 ])
7 df_location_points = pd.read_csv('/content/gdrive/My Drive/A_Doctorado/datos_thesis/Pleiades_CSVs/location_points.csv', usecols=[
8 | 'place_id', 'geometry_wkt'
9 ])
10 df_location_polygons = pd.read_csv('/content/gdrive/My Drive/A_Doctorado/datos_thesis/Pleiades_CSVs/location_polygons.csv', usecols=[
11 | 'place_id', 'geometry_wkt'
12 ])
13
14 # Filtro para encontrar las entidades nombradas
15 entidades = ['Castalia', 'Scythiae', 'Veios', 'Iasonia', 'Assyrios', 'Rhamnusia', 'Lethaeam', 'Romae', 'Spartana', 'Castra', 'Asis', 'Creta', 'Tyrrhenos', 'Calliope', 'Thraciam', '
16
```

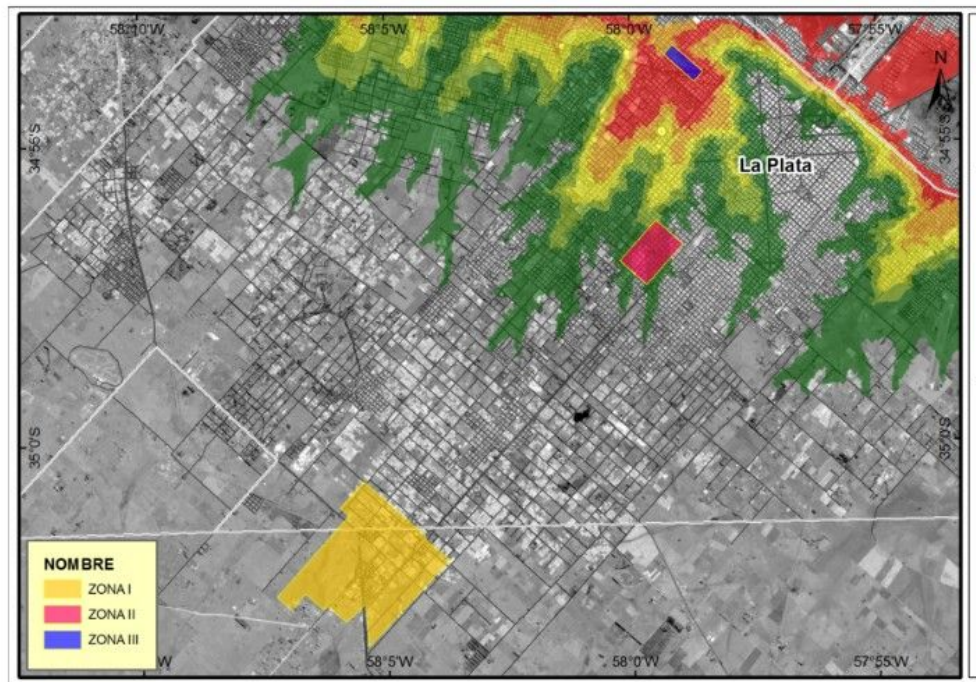
Processing Data from Pleiades

	uri	place_id	attested_form	\
0	https://pleiades.stoa.org/places/993/gallia	993	NaN	
1	https://pleiades.stoa.org/places/991401/august...	991401	NaN	
2	https://pleiades.stoa.org/places/991409/august...	991409	NaN	
3	https://pleiades.stoa.org/places/991373/creta	991373	NaN	
4	https://pleiades.stoa.org/places/991386/phrygia-i	991386	NaN	

	romanized_form_1	romanized_form_2	romanized_form_3	entidad_buscada	\
0	Gallia	NaN	NaN	Gallia	
1	Augusta Euphratensis	NaN	NaN	Euphraten	
2	Augusta Libanensis	NaN	NaN	Liba	
3	Creta	NaN	NaN	Creta	
4	Phrygia I	NaN	NaN	Phrygia	

	geometry_wkt	geometry_wkt_polygon
0	NaN	NaN
1	NaN	POLYGON ((35 35, 40 35, 40 40, 35 40, 35 35))
2	NaN	POLYGON ((35 30, 40 30, 40 35, 35 35, 35 30))
3	NaN	POLYGON ((20 35, 25 35, 25 40, 20 40, 20 35))
4	NaN	POLYGON ((25 35, 30 35, 30 40, 25 40, 25 35))

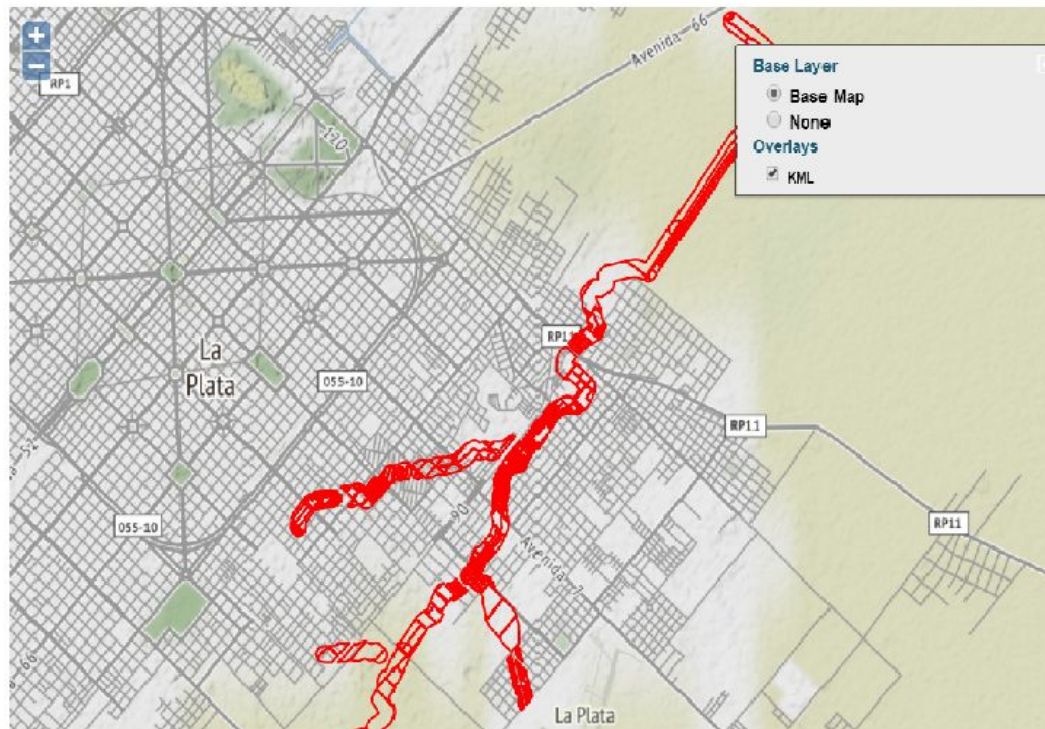
Processing Data from Pleiades



Processing Data from Pleiades



Processing Data from Pleiades



Wikidata API



- Wikidata is a free, collaborative, multilingual database that serves as a **central repository for structured data**.
- Wikidata offers open, structured data that can be linked to various projects, making it ideal for research, data analysis, and automation.
- It connects data from other Wikimedia projects (Wikipedia, Wikimedia Commons) and external databases, serving as a hub for cross-referencing.
- Wikidata **provides unique identifiers** (Q numbers) for **people, places, objects, and concepts**, allowing easy access and integration across datasets.

Wikidata API



```
3 import requests
4
5 def search_wikidata(place_name):
6     url = "https://www.wikidata.org/w/api.php"
7     params = {
8         "action": "wbsearchentities",
9         "language": "en",
10        "format": "json",
11        "search": place_name,
12        "type": "item"
13    }
14    try:
15        response = requests.get(url, params=params)
16        data = response.json()
17        results = []
18        if data['search']:
19            for result in data['search']:
20                details = get_entity_images(result['id'])
21                results.append({
22                    "id": result['id'],
23                    "label": result['label'],
24                    "images": details
25                })
26            return results
27        else:
28            return []
29    except requests.RequestException as e:
30        return {"error": f"Request failed: {e}"}
31
```

Wikidata API



```
32 def get_entity_images(entity_id):
33     url = "https://www.wikidata.org/w/api.php"
34     params = {
35         "action": "wbgetclaims",
36         "entity": entity_id,
37         "property": "P18",
38         "format": "json"
39     }
40     try:
41         response = requests.get(url, params=params)
42         data = response.json()
43         images = []
44         if 'claims' in data and 'P18' in data['claims']:
45             for claim in data['claims']['P18'][:10]: # Limit to 10 images
46                 image_title = claim['mainsnak']['datavalue']['value']
47                 image_url = f"https://commons.wikimedia.org/wiki/File:{image_title.replace(' ', '_')}}"
48                 images.append(image_url)
49         return images
50     except requests.RequestException:
51         return []
```

Wikidata API

Results for Scythiae:

Results for Veios:

- Veii (Q677935):

https://commons.wikimedia.org/wiki/File:Colonne_di_Veio_a_palazzo_Wedekind.JPG

<https://commons.wikimedia.org/wiki/File:Veii.png>

Results for Iasonia:

- Jasonia (Pontus) (Q17635213):

Results for Assyrios:

Results for Rhamnusia:

- Rhamnusia Virgo (Q59686630):

- Nemesis (Q185747):

https://commons.wikimedia.org/wiki/File:Statue_Nemesis_Louvre_Ma4873.jpg

Results for Lethaeam:

Results for Romae:

- Parilia (Q666621):

https://commons.wikimedia.org/wiki/File:Suvé_e_Festa_di_Pales.jpg

- Romance Is a Bonus Book (Q60174208):

- Romaeuropa Festival (Q55829521):

- Romaero (Q387165):

https://commons.wikimedia.org/wiki/File:Baneasa_Airport_4.jpg

- Romaen Campbell (Q116476699):

- Romani people (Q8060):

https://commons.wikimedia.org/wiki/File:Bain_rituel_dans_la_mer_aux_Saintes-Maries.jpg

[https://commons.wikimedia.org/wiki/File:Romanlar_\(Istanbul\).JPG](https://commons.wikimedia.org/wiki/File:Romanlar_(Istanbul).JPG)

- Romeo and Juliet (Q83186):

https://commons.wikimedia.org/wiki/File:Romeo_and_juliet_brown.jpg

https://commons.wikimedia.org/wiki/File:Romeo_Juliet.jpg

<https://commons.wikimedia.org/wiki/File:DickseeRomeoandJuliet.jpg>

<https://commons.wikimedia.org/wiki/File:Rigaud-RomeoJuliet.jpg>

Results for Spartana:

Wikidata API

Results for Scythiae:

Results for Veios:

- Veii (Q677935):

https://commons.wikimedia.org/wiki/File:Colonne_di_Veio_a_palazzo_Wedekind.JPG

<https://commons.wikimedia.org/wiki/File:Veii.png>

Results for Iasonia:

- Jasonia (Pontus) (Q17635213):

Results for Assyrios:

Results for Rhamnusia:

- Rhamnusia Virgo (Q59686630):

- Nemesis (Q185747):

https://commons.wikimedia.org/wiki/File:Statue_Nemesis_Louvre_Ma4873.jpg

Results for Lethaeam:

Results for Romae:

- Parilia (Q666621):

https://commons.wikimedia.org/wiki/File:Suvé_e_Festa_di_Pales.jpg

- Romance Is a Bonus Book (Q60174208):

- Romaeuropa Festival (Q55829521):

- Romaero (Q387165):

https://commons.wikimedia.org/wiki/File:Baneasa_Airport_4.jpg

- Romaen Campbell (Q116476699):

- Romani people (Q8060):

https://commons.wikimedia.org/wiki/File:Bain_rituel_dans_la_mer_aux_Saintes-Maries.jpg

[https://commons.wikimedia.org/wiki/File:Romanlar_\(Istanbul\).JPG](https://commons.wikimedia.org/wiki/File:Romanlar_(Istanbul).JPG)

- Romeo and Juliet (Q83186):

https://commons.wikimedia.org/wiki/File:Romeo_and_juliet_brown.jpg

https://commons.wikimedia.org/wiki/File:Romeo_Juliet.jpg

<https://commons.wikimedia.org/wiki/File:DickseeRomeoandJuliet.jpg>

<https://commons.wikimedia.org/wiki/File:Rigaud-RomeoJuliet.jpg>

Results for Spartana:

Wikidata API

File:Statue Nemesis Louvre Ma4873.jpg

From Wikimedia Commons, the free media repository

[Download](#) [Use this file](#) [Use this file](#) [Email a link](#) [Information](#)



The moral of this experiment

- Some tasks, like the automatic generation of TEI files and automatic tagging, work very well and streamline the editing process.
- However, they still require a **final review by the editor** to ensure accuracy and quality control, making the human touch essential for the last stage of refinement.
- The three prototypes for improving the digital edition were successful, as **we successfully** integrated **APIs** from **VIAF**, **Pleiades**, and **Wikidata**.
- However, we have become victims of our own success. Now, we face a **new bottleneck**: processing **the massive amount of data we've retrieved**.
- With thousands of names, places, and images, **the challenge** is no longer in retrieving or manually searching for data but in **efficiently processing and filtering the correct information locally**.
- Managing this overwhelming volume of data requires new strategies and tools for optimization.



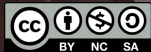
UNIVERSIDAD
NACIONAL
DE LA PLATA

TEI 2024

Carlos J. Nusch

carlosnusch@prebi.unlp.edu.ar

PREBI SEDICI
prebi.unlp.edu.ar sedici.unlp.edu.ar



Esta obra está bajo una Licencia Creative Commons
Atribución-NoComercial-CompartirIgual 4.0 Internacional

TEI

Text
Encoding
Initiative



USAL
UNIVERSIDAD
DEL SALVADOR

**THANK
YOU**

