

# Encyclopædia Britannica;

James OR, A Fullerton

## DICTIONARY

O F

A R T S and S C I E N C E S,

COMPILED UPON A NEW PLAN.

I N W H I C H

The diferent SCIENCES and ARTS are digested into  
distinct Treatises or Systems;

A N D

The various TECHNICAL TERMS, &c. are explained as they occur  
in the order of the Alphabet.

ILLUSTRATED WITH ONE HUNDRED AND SIXTY COPPER

By a SOCIETY of GENTLEMEN in SCOTLA

I N T H R E E V O L U M E S

V O L. I.

E D I N B U R

Printed for A. BELL and C. N.  
And sold by COLIN MACFARQUHAR, at

M. DCC

***Frances: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale***

**Dr. Rosa Filgueira,  
Lecturer at the School of Computer Science,  
University of St Andrews,  
Email: rf208@st-Andrews.ac.uk**

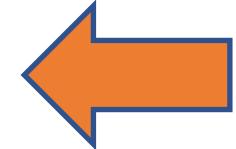
# Background – Research Interest

## Scalability, reliability, and performance of data science applications

- Data/task workflows & programming abstractions
- Automatic and portable parallelisation techniques
- Workload strategies that automatically adapt to applications at runtime
- Dynamic optimization techniques for stream-based workflows/frameworks
- Techniques for dealing with failures at the application and the system levels
- Exploiting heterogeneous platforms for data-intensive applications

## Unlocking Historical Digital Text Collections through Advanced AI methods and Parallel Techniques

- Large scale text mining facilities in digital humanities
- Hiding system complexity
- Ease of use and consistency for complex text-mining analysis
- Techniques for making digital text searchable, and analysable
- Complex text analysis queries automatically and at scale, and for visualizing results
- Advanced AI-techniques to extract knowledge from digital text collections.



## Intelligent advisory systems for scientific software

- Software feature extraction
- Static code analysis techniques
- Semantic code search algorithms; automatic code classification/tagging methods;
- Automatic pipeline/workflow composition and orchestration;
- -Automated serverless computing architectures

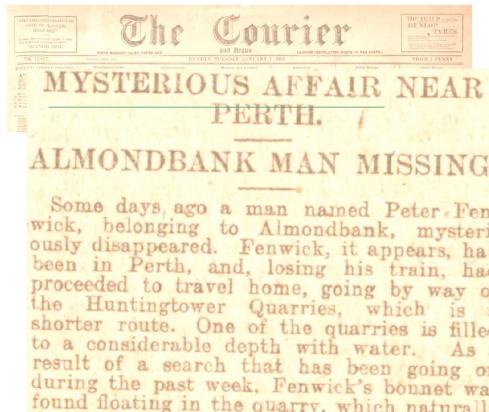
# frances: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale

## Motivation

- Large digital collections been available for research
- Hunger for large scale text mining facilities and for:
  - HPC/Cloud environments
  - Analytic frameworks to create applications

## Challenges

- Several large digital collections (semi-structured data)
- Different levels of quality of data – OCR
- Data with different physical representations and XML schemas
- Difficult to mine, search, analyse, compare



## Collaborations ~ 6 years:

- *Melisa Terras, Lisa Otty* – CDCS UoE
- *Bea Alex, Sarah van Eindhoven, Lisa Gotthard* UoE
- *Sarah Ames, NLS*
- *LwM project- Alan Turing*

```
...
<text.title>
  <pg pgref="5" clipref="1"
       pos="4069,3036,4949,3154"/>
  <p>
    <wd pos="4069,3036,4949,3154">MYSTERIOUS AFFAIR
    NEAR PERTH.</wd>
  </p>
</text.title>
<text.cr>
  <pg pgref="5" clipref="1"
       pos="4039,3191,4987,4235"/>
  <p>
    <wd pos="4041,3192,4496,3241">ALMONDBANK</wd>
    <wd pos="4523,3200,4663,3246">MAN</wd>
    <wd pos="4696,3198,4976,3250">MISSING.</wd>
    <wd pos="4085,3290,4189,3323">Some</wd>
    <wd pos="4214,3290,4312,3329">days,</wd>
  ...

```

# *frances*: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale



[Using the Library](#) | [Digital resources](#) | [Catalogues](#) | [Shop](#) | [News & events](#) | [Contact](#)

[Join / Log in](#)



[Home](#) › [Using the Library](#) › [Academic research](#) › [Fellowships](#) › [Digital scholarship](#)

## The National Librarian's Research Fellowship in Digital Scholarship

Fellowships

Graham Brown

Digital scholarship

### 2021-2022 National Library of Scotland Digital Scholarship Fellow

Dr Rosa Filgueira is an Assistant Professor with the School of Mathematical and Computer Sciences at Heriot-Watt University.

Her research background includes high-performance computing, data streaming, data-intensive computing, and large-scale distributed systems.

Rosa has experience handling and mining large digital textual collections. She has worked on several digital humanities and data science projects, co-developing and co-designing text-mining applications at scale.



Dr Rose Filgueira

#### AI toolbox project

This project will explore new ways to unlock the full value of the National Library of Scotland's [Data Foundry](#) collections by building a new AI toolbox called 'frances' and a web user interface that allows researchers to interact with it.

It would help historians (and other users) to extract complex knowledge from digital collections in a fast and transparent manner without having to be an expert data scientist. The toolbox is named after Frances Wright (September 6, 1795 – December 13, 1852), a Scottish-born lecturer, writer,

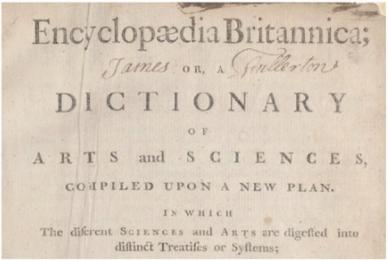
\*[Frances Wright](#) (September 6, 1795 – December 13, 1852), was a Scottish-born lecturer, writer, freethinker, feminist, utopian socialist, abolitionist, social reformer, and Epicurean philosopher

# frances: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale

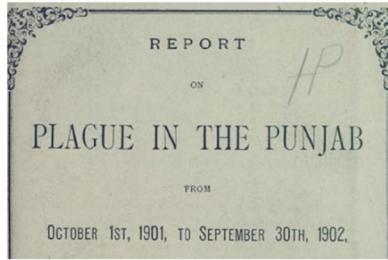
## Digitised collections

Download the ALTO, METS, image and plain text files for our digitised collections.

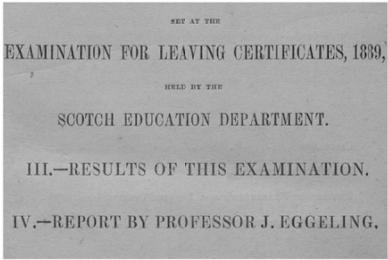
**8 Editions:**  
1768 -1860



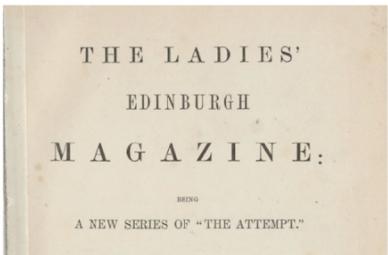
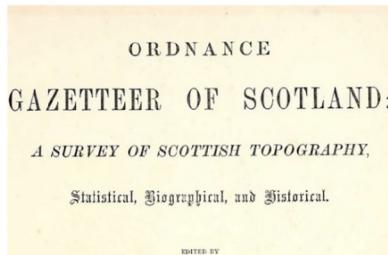
Encyclopaedia Britannica, 1768-1860



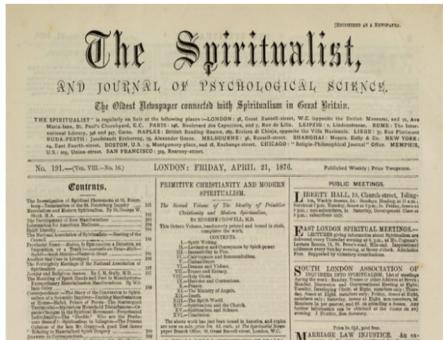
A Medical History of British India



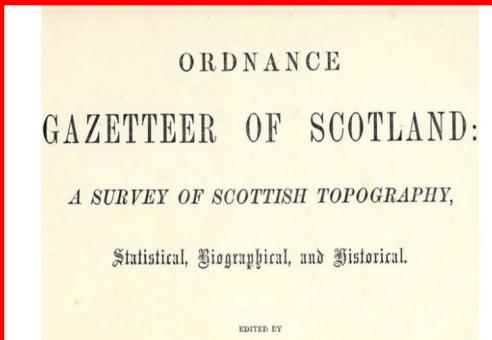
Scottish School Exam Papers, 1888-1963



# frances: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale



The Spiritualist



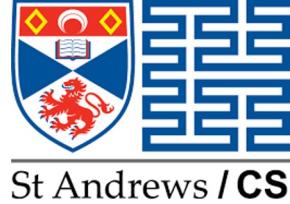
Asian Directories and Chronicles



Britain and UK Handbooks



# *frances: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale*



Goal: Provide **abstractions** to a variety of ML/NLP techniques

→ Extract complex knowledge without being an expert data scientist

- Train and use text embedding models
- Employ topic mining, sentiment analysis, text summarization
- Build knowledge graph(s) visualizing the results
- Text Mining Parallel Platform (based in Apache Spark)
- Automatic Visualization of Results



New suite of ML functionalities that can be used to analyse  
any other Data Foundry collection

## The Journey

Phase 1: Information Extraction

Phase 2: Ontologies and Knowledge Graphs

Phase 3: Augmented Knowledge Graphs with Deep Transfer Learning

Phase 4: Defoe and Knowledge Graphs

Phase 5: React-Flask Web Platform

Phase 6: Case Studies

## The Journey

### Phase 1: Information Extraction

#### 1.1 Encyclopedia Britannica

#### 1.2 Others

Phase 2: Ontology and Knowledge Graphs

Phase 3: Augmented Knowledge Graphs with Deep Transfer Learning

Phase 4: Defoe and Knowledge Graphs

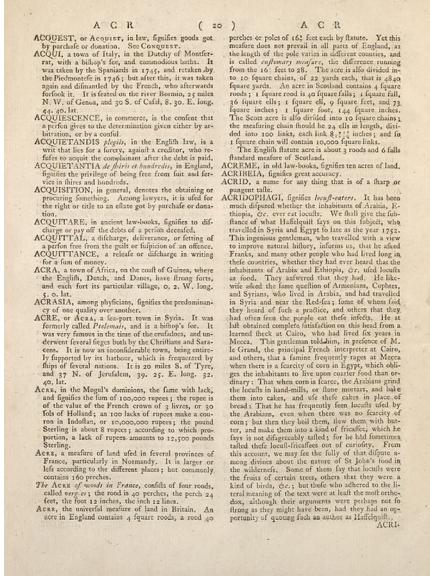
Phase 5: React-Flask Web Platform

Phase 6: Case Studies

## Detect, Classify and Extract all EB terms across editions

- Two types of Terms:
  - **Articles:** 1 or 2 paragraphs describing a term ~ entry in a dictionary
  - **Topics:** Several pages describing a term
- Extract Metadata from all EB editions and volumes.

How have we done it ? By using *defoe* to analyze METS and ALTO XML files



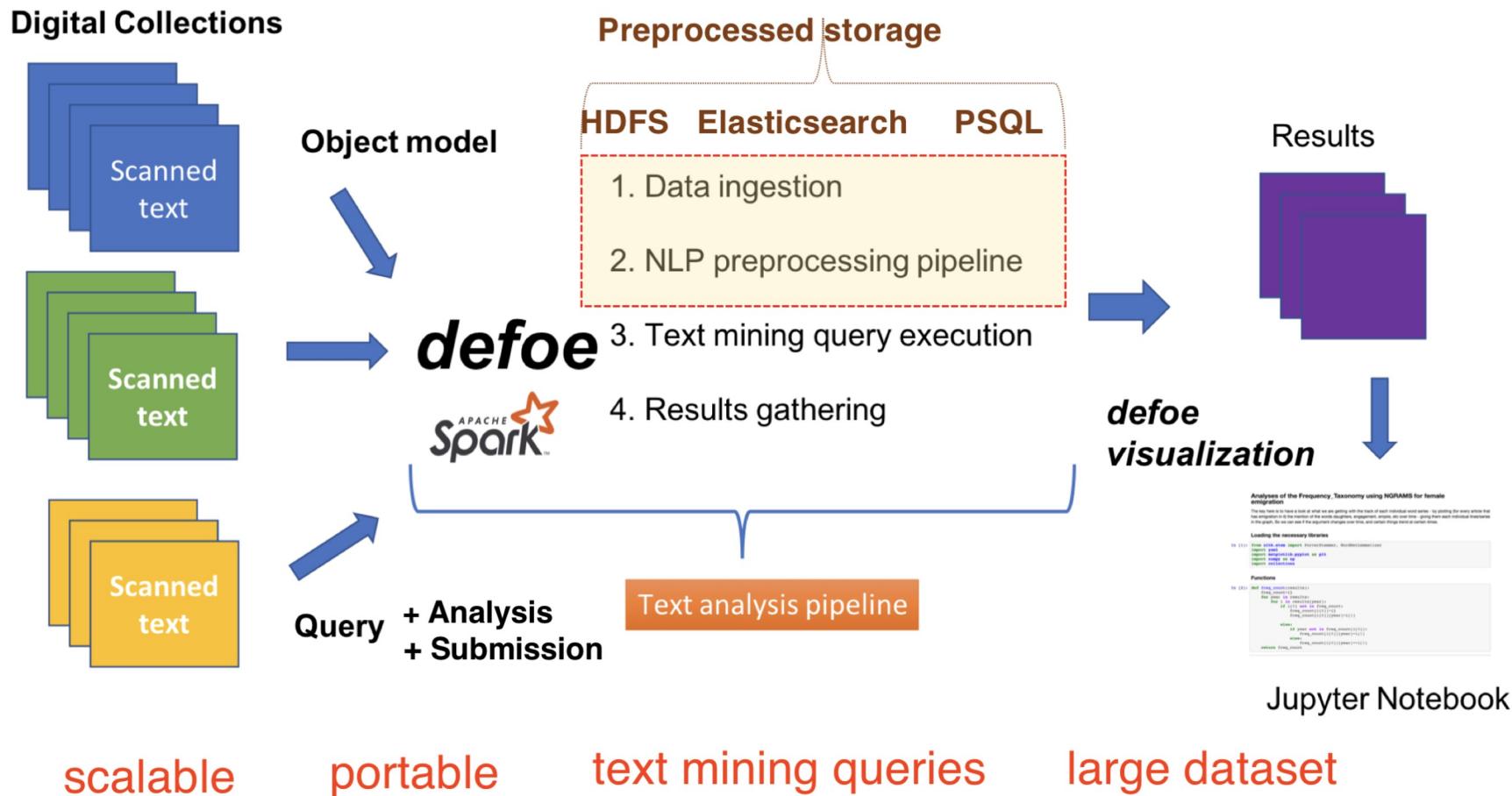
Page 36 of the First Edition  
of the Encyclopaedia Britannica

```

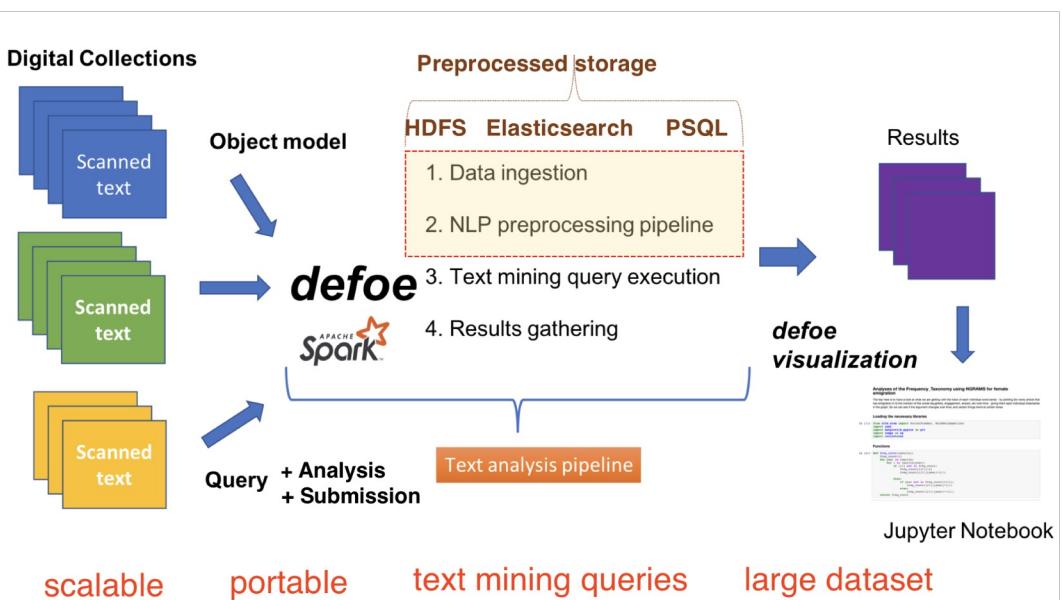
1 <?xml version="1.0" encoding="UTF-8"?>
2 <alto
3   xmlns="http://www.loc.gov/standards/alto/v3/alto.xsd">
4   <@descript>
5     <@documentUnit:pixel></MeasurementUnit>
6     <@sourceImageInformation>
7       <@fileName>/data/dfs/c_188936619/l_141133901/18882693.23.pdf</fileName>
8     </@sourceImageInformation>
9     <@ocrProcessingStep ID="1000">
10    <@processingDate>Fri Jul 12 10:36:14 2019
11    <@processingSoftware>
12      <@softwareCreator>CONTRIBUTORS</softwareCreator>
13      <@softwareName>pdftoalto</softwareName>
14      <@softwareVersion>1.0.0</softwareVersion>
15    </@processingSoftware>
16  </@ocrProcessingStep>
17  </@descript>
18 </alto>
19 <@processingStep>
20   <@processingSoftware>
21     <@softwareCreator>CONTRIBUTORS</softwareCreator>
22     <@softwareName>pdftoalto</softwareName>
23     <@softwareVersion>1.0.0</softwareVersion>
24   </@processingSoftware>
25 </@processingStep>
26 <@Layout>
27   <@Page ID="Page36" PHYSICAL_IMG_NR="36" WIDTH="2466" HEIGHT="3327">
28     <@TextLine WIDTH="232" HEIGHT="75" ID="p36_l1" HPOS="512" VPOS="180">
29       <String ID="p36_w1" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
30     <@TextLine WIDTH="57" VPOS="187" HPOS="554"/>
31     <String ID="p36_w2" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
32     <@TextLine WIDTH="182" VPOS="187" HPOS="648"/>
33     <String ID="p36_w3" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
34     <@TextLine WIDTH="132" HEIGHT="75" ID="p36_l2" HPOS="1119" VPOS="180">
35       <String ID="p36_w4" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
36     <@TextLine WIDTH="62" VPOS="187" HPOS="1138"/>
37     <String ID="p36_w5" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
38     <@TextLine WIDTH="182" HEIGHT="75" ID="p36_l3" HPOS="171" VPOS="182">
39       <String ID="p36_w6" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
40     <@TextLine WIDTH="243" HEIGHT="75" ID="p36_l4" HPOS="182" VPOS="182">
41       <String ID="p36_w7" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
42     <@TextLine WIDTH="182" HEIGHT="75" ID="p36_l5" HPOS="414" VPOS="182">
43       <String ID="p36_w8" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
44     <@TextLine WIDTH="15" VPOS="182" HPOS="477"/>
45       <String ID="p36_w9" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
46     <@TextLine WIDTH="182" HEIGHT="75" ID="p36_l6" HPOS="492" VPOS="182">
47       <String ID="p36_w10" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
48     <@TextLine WIDTH="124" HEIGHT="75" ID="p36_l7" HPOS="685" VPOS="182">
49       <String ID="p36_w11" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
50     <@TextLine WIDTH="134" HEIGHT="75" ID="p36_l8" HPOS="831" VPOS="182">
51       <String ID="p36_w12" CONTENT="Acquisition, or Acquirer, in law, signifies goods got by purchase or donation, in the exercise of a right, or by force; as, the property of the victor in war, or of the conqueror in 1765; but after this, it was taken up, and different meanings given to it, according to the sense in which it is used. It is granted on the river Rhenus, 21 miles N. W. of Geneva, and 25 of Cufa, 8. 30. E. long. 46. 40. N. lat."/>
52     <@TextLine WIDTH="175" HEIGHT="75" ID="p36_l9" HPOS="991" VPOS="182">
53       <String ID="p36_w13" CONTENT="got," HPOS="1117" VPOS="182"/>
54     <@TextLine WIDTH="59" HEIGHT="75" ID="p36_l10" HPOS="1141" VPOS="182"/>
55   </TextLine>

```

ALTO-XML of the Page 36 of the First Edition of the Encyclopaedia Britannica



<https://github.com/francesNLP/defoe>



## Publications:

1. [defoe: A Spark-based Toolbox for Analysing Digital Historical Textual Data](#), 2019
2. [Geoparsing the historical gazetteers of scotland: Accurately computing location in mass digitised texts](#), 2020
3. [Extending defoe for the efficient analysis of historical texts at scale](#), 2021

## New defoe queries and heuristics

- **defoe query to extract EB Terms:**
  - By page and classify them between **articles** and **topics** → Using ALTO XML

A B A ( 2 ) A B B

fome orhamant, as a rose or other flower. Scamozzi uses *abacis* for a concave ornament on the capital of an Italian pedestal. In heraldry, the epithet applied to the wings of eagles, &c., when the top hole of each is to the right of the field, or when the wings are flat; the natural way of bearing them being extended.

ABASING, in the fra-language, signifies the same as striking.

ABASSI, or *ABAAS*, a silver coin current in Persia, equivalent in value to a French livre, or a pence half-penny. Struck by Shah Abbas II, King of Persia, under whom it was struck.

ABATAMENTUM, in law, is an entry to lands by interposition, i.e., when a person dies seized, and another who has no right enters before the heir.

ABATE, from *abatre*, to destroy; a term used by writers on the canon law, both in an active and neutral sense. To *abate* a calle, is to destroy or beat it down; to *abate* a wint, is, by some exception to render it null and void.

ABATE, in the manage, implies the performance of any downward motion properly. Hence a horse is said to *abate*, or take down his curves, when he puts both his fore legs to the ground at once, and observes the same exercise in all the times.

ABATEMENT, in heraldry, implies something added to a coat of arms in order to lessen its dignity, and point out some imperfection or flaw in the character of the wearer.

ABATEMENT, in law. See *ABATE*.

ABATEMENT, in commerce, signifies an allowance or discount in the price of certain commodities, in consideration of prompt payment; a diminution in the fluctuating quantity or quality of goods, or some such circumstance.

ABATEMENT, in the customs, an allowance made upon the duty of goods, when the quantum damaged is determined by the judgment of two merchants upon oath, and ascertained by a certificate from the surveyor and waiver.

ABATIS, an ancient term for an officer of the stables.

ABATOR, in law, a term applied to a person who enters to a house or lands, void by the death of the last possessor, before the true heir.

ABAVO, in botany, a synonyme of the adansonia, a genus of plants belonging to the icacandia polygalia of Linnaeus. See CALYCANTHUS.

ABAISSE. See *ABASE*.

ABALIENATION. See *ALIENATION*.

ABANBO, a river of Ethiopia which falls into the Nile.

ABANCAL, or ABANCAYA, a town and river of Peru, in the district of Lima.

ABANCO, a town of Italy, subject to Venice, and situated five miles south-west of Padua.

ABAPTISTON, or ANABAPTISTON, an obsolete term for the chirurgical instrument called a *tropan*. See *Suspensory*, and *Trepan*.

ABARCA, a shoe made of raw hides, formerly worn by the peasants in Spain.

ABARTICULATION, in anatomy, a species of articulation, so termed *duribracta*. See *Anatomy*, Part I, and *Diseases*.

ABAS, a weight used in Persia for weighing pearls. It is 1-8th less than the European carat.

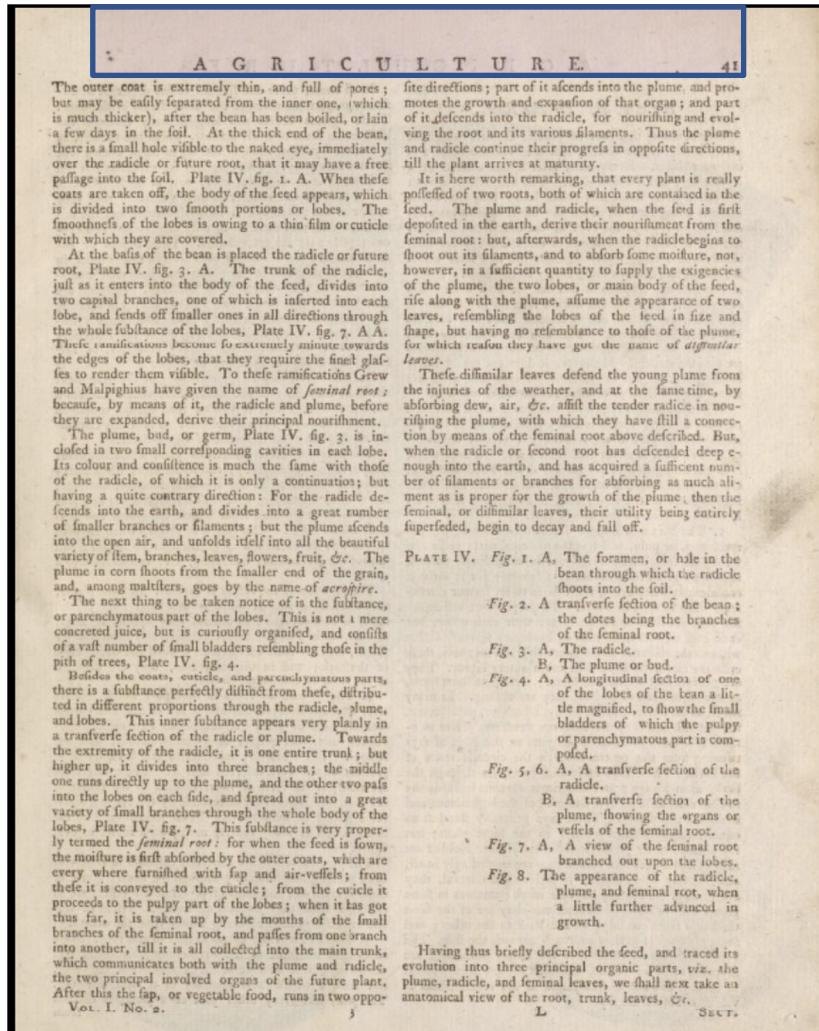
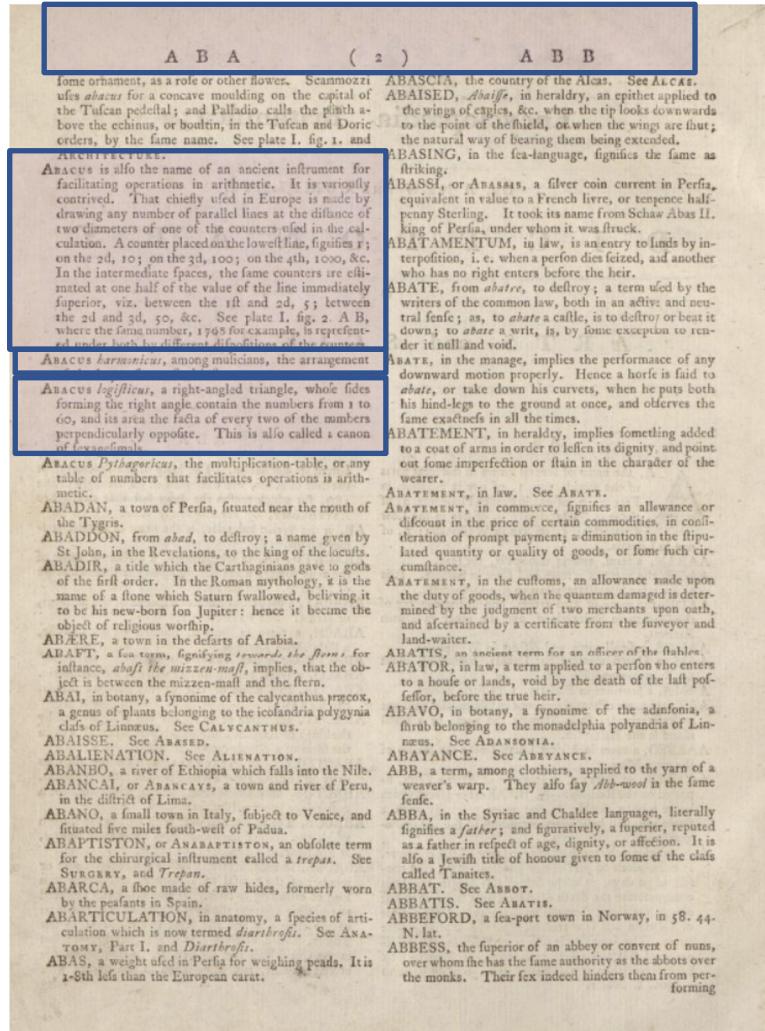
```

    ↳ Layout
      ↳ Page
        ↳ PrintSpace
          ↳ TextLine
            ↳ String : (@ID = p16_w1,@CONTENT = ABA,@HPOS = 509,@VPOS = 160,@WIDTH = 233,@HEIGHT = 45,@STYLEREFS = font0)
        ↳ TextLine
          ↳ String : (@ID = p16_w2,@CONTENT = fome,@HPOS = 175,@VPOS = 238,@WIDTH = 84,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 13,@VPOS = 238,@HPOS = 259)
          ↳ String : (@ID = p16_w3,@CONTENT = ornament,,@HPOS = 272,@VPOS = 238,@WIDTH = 178,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 15,@VPOS = 238,@HPOS = 450)
          ↳ String : (@ID = p16_w4,@CONTENT = as,@HPOS = 465,@VPOS = 238,@WIDTH = 34,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 13,@VPOS = 238,@HPOS = 499)
          ↳ String : (@ID = p16_w5,@CONTENT = a,@HPOS = 512,@VPOS = 238,@WIDTH = 19,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 10,@VPOS = 238,@HPOS = 531)
          ↳ String : (@ID = p16_w6,@CONTENT = rofe,@HPOS = 541,@VPOS = 238,@WIDTH = 69,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 238,@HPOS = 610)
          ↳ String : (@ID = p16_w7,@CONTENT = or,@HPOS = 623,@VPOS = 238,@WIDTH = 36,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 12,@VPOS = 238,@HPOS = 659)
          ↳ String : (@ID = p16_w8,@CONTENT = other,@HPOS = 671,@VPOS = 238,@WIDTH = 93,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 11,@VPOS = 238,@HPOS = 764)
          ↳ String : (@ID = p16_w9,@CONTENT = flower,@HPOS = 775,@VPOS = 238,@WIDTH = 122,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 50,@VPOS = 238,@HPOS = 897)
          ↳ String : (@ID = p16_w10,@CONTENT = Scamozzi,@HPOS = 947,@VPOS = 238,@WIDTH = 207,@HEIGHT = 37,@STYLEREFS = font1)
        ↳ TextLine
          ↳ String : (@ID = p16_w11,@CONTENT = ufes,@HPOS = 174,@VPOS = 284,@WIDTH = 69,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 16,@VPOS = 284,@HPOS = 243)
          ↳ String : (@ID = p16_w12,@CONTENT = abacus,@HPOS = 259,@VPOS = 284,@WIDTH = 116,@HEIGHT = 37,@STYLEREFS = font2)
          ↳ SP : (@WIDTH = 16,@VPOS = 284,@HPOS = 375)
          ↳ String : (@ID = p16_w13,@CONTENT = for,@HPOS = 391,@VPOS = 284,@WIDTH = 53,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 16,@VPOS = 284,@HPOS = 444)
          ↳ String : (@ID = p16_w14,@CONTENT = a,@HPOS = 460,@VPOS = 284,@WIDTH = 18,@HEIGHT = 37,@STYLEREFS = font1)
          ↳ SP : (@WIDTH = 16,@VPOS = 284,@HPOS = 478)
          ↳ String : (@ID = p16_w15,@CONTENT = concave,@HPOS = 494,@VPOS = 284,@WIDTH = 137,@HEIGHT = 37,@STYLEREFS = font1)
        ↳ SP : (@WIDTH = 22,@VPOS = 284,@HPOS = 521)
    
```

# Phase 1.1: Information Extraction - EB

## New defoe queries and heuristics

1. Detecting pages headers from ALTO XML
2. Using headers to classify terms into: Articles & Topics
3. Using ALTO Text for detecting the start of each article:  
--> Starting a line with TERM UPPERCASE + “,”



Articles

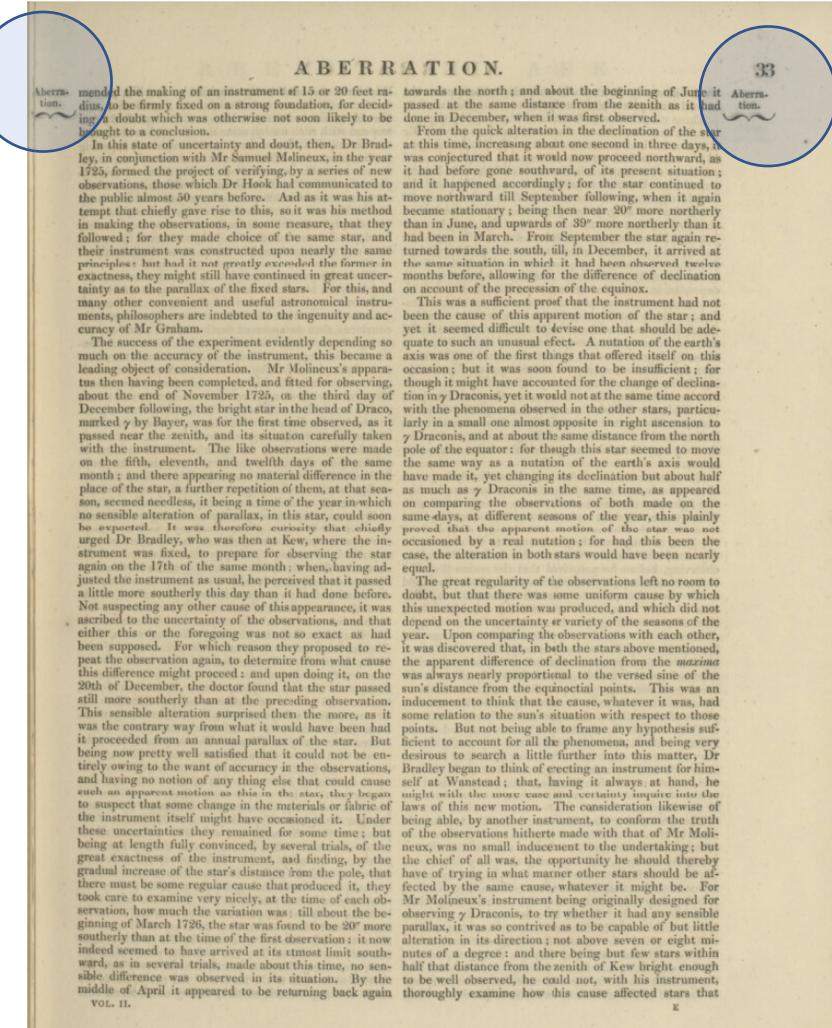
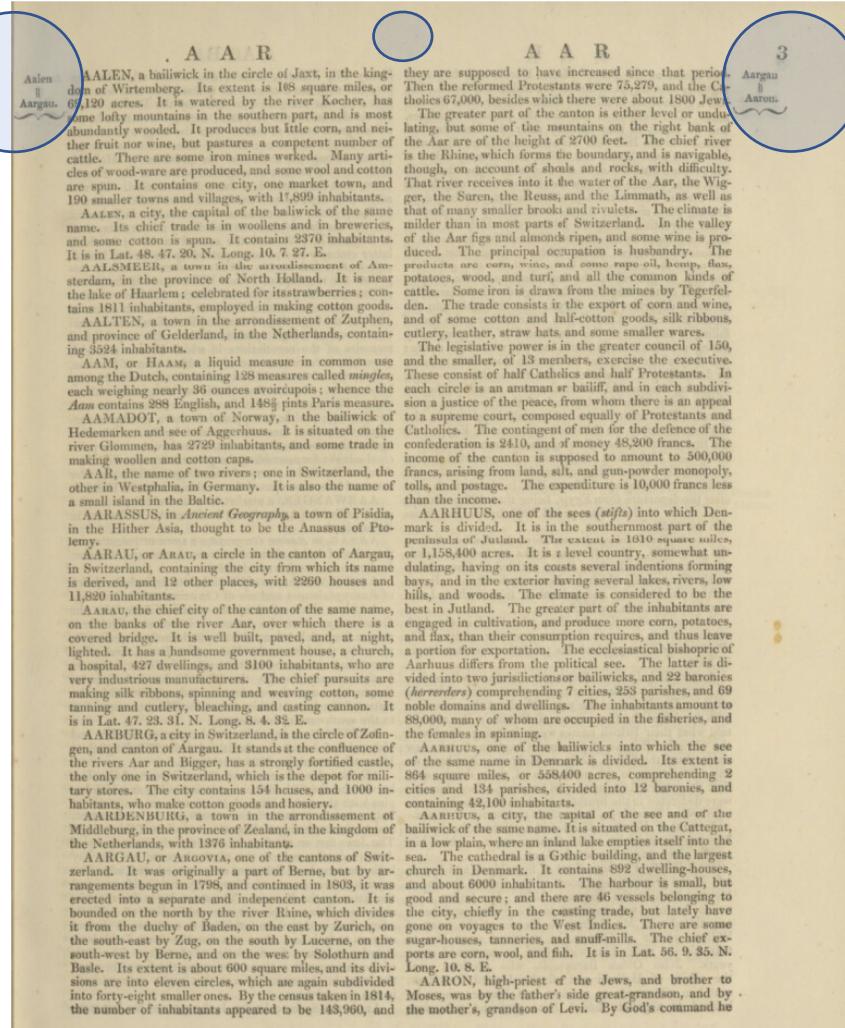
Edition 1 - 1771

Topic

# Phase 1.1: Information Extraction - EB

## New defoe queries and heuristics

1. Detecting pages headers from ALTO XML
2. Using headers to classify terms into: Articles & Topics
3. Using ALTO Text for detecting the start of each article:  
--> Starting a line with TERM UPPERCASE + “,”



Articles

Edition 7

Topic

## New *defoe* queries and heuristics

term		SCIENCE
definition	in philosophy, denotes any ddpfrine, deduced f...	
relatedTerms		[]
header		SCISCO
startsAt		658
endsAt		658
numberOfTerms		24
numberOfWords		15
numberOfPages		872
positionPage		7
typeTerm		Article
editionTitle	First edition, 1771, Volume 3, M-Z	
editionNum		1
supplementTitle		
supplementsTo		[]
year		1771
place		Edinburgh
volumeTitle	Encyclopaedia Britannica; or, A dictionary of ...	
volumeNum		3
letters		M-Z
part		0
altoXML	144133903/alto/144812443.34.xml	
Name: 7454, dtype: object		

**Science term** information extracted from the First Edition - 1771

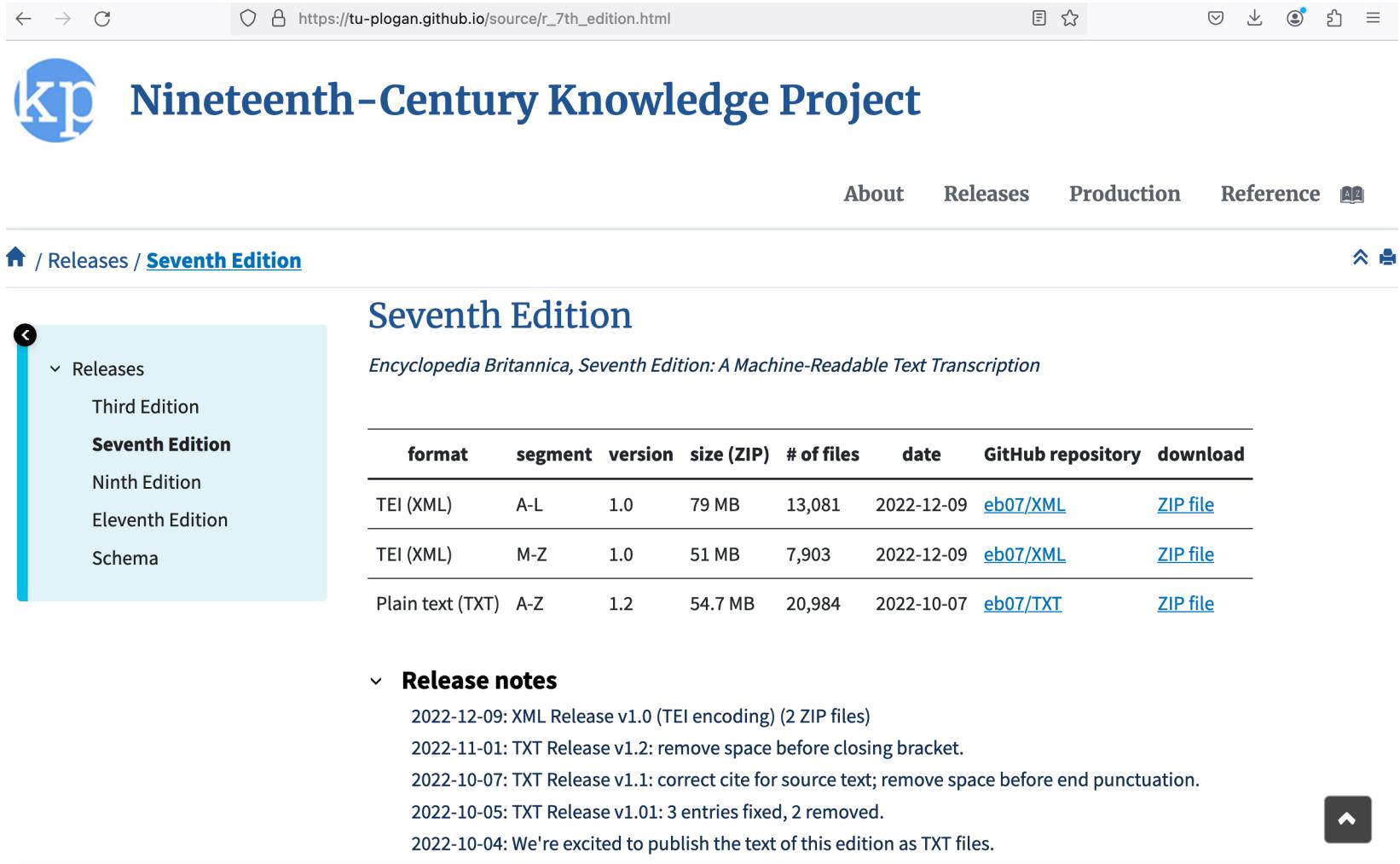
## New *defoe* queries and heuristics

- *defoe* query to extract EB metadata :
  - Metadata per Edition and Volume → Using METS XML

	MMSID	editionTitle	editor	editor_date	genre	language	termsOfAddress	numberOfPages	physicalDescription	place	...
0	992277653804341	First edition, 1771, Volume 1, A-B	Smellie, William	1740-1795	encyclopedia	eng	None	832	3 v., 160 plates : ill. ; 26 cm. (4to)	Edinburgh	... h
1	992277653804341	First edition, 1771, Volume 2, C-L	Smellie, William	1740-1795	encyclopedia	eng	None	1018	3 v., 160 plates : ill. ; 26 cm. (4to)	Edinburgh	... h
2	992277653804341	First edition, 1771, Volume 3, M-Z	Smellie, William	1740-1795	encyclopedia	eng	None	872	3 v., 160 plates : ill. ; 26 cm. (4to)	Edinburgh	... h
3	9929192893804340	First edition, 1773, Volume 1, A-B	Smellie, William	1740-1795	encyclopedia	eng	None	844	3 v. (viii, 697, [1] p., LVIII leaves of plate...)	London	... h
4	9929192893804340	First edition, 1773, Volume 2, C-L	Smellie, William	1740-1795	encyclopedia	eng	None	1032	3 v. (viii, 697, [1] p., LVIII leaves of plate...)	London	... h
5	9929192893804340	First edition, 1773, Volume 3, M-Z	Smellie, William	1740-1795	encyclopedia	eng	None	864	3 v. (viii, 697, [1] p., LVIII leaves of plate...)	London	... h

Subset of **metadata extracted** for the volumes of the **First Edition**. This edition is a 3-volume reference work, **issued twice**, in **1771 and 1773**.

## Work in progress- Using the clean data from the 7<sup>th</sup> Edition - KP



The screenshot shows a web browser displaying the [Nineteenth-Century Knowledge Project](https://tu-plogen.github.io/source/r_7th_edition.html) website. The URL in the address bar is [https://tu-plogen.github.io/source/r\\_7th\\_edition.html](https://tu-plogen.github.io/source/r_7th_edition.html). The page title is "Seventh Edition". The navigation menu includes links for About, Releases, Production, Reference, and a sorting icon. The breadcrumb navigation shows the user is at the home page, then Releases, and finally Seventh Edition. On the left, a sidebar menu is open, showing a list of releases: Releases, Third Edition, **Seventh Edition**, Ninth Edition, Eleventh Edition, and Schema. The main content area displays the "Seventh Edition" page, which includes a subtitle "Encyclopedia Britannica, Seventh Edition: A Machine-Readable Text Transcription". Below this is a table of contents and a table of file downloads. The table has columns for format, segment, version, size (ZIP), # of files, date, GitHub repository, and download link. The data is as follows:

format	segment	version	size (ZIP)	# of files	date	GitHub repository	download
TEI (XML)	A-L	1.0	79 MB	13,081	2022-12-09	<a href="#">eb07/XML</a>	<a href="#">ZIP file</a>
TEI (XML)	M-Z	1.0	51 MB	7,903	2022-12-09	<a href="#">eb07/XML</a>	<a href="#">ZIP file</a>
Plain text (TXT)	A-Z	1.2	54.7 MB	20,984	2022-10-07	<a href="#">eb07/TXT</a>	<a href="#">ZIP file</a>

Below the table, under the heading "Release notes", are the following entries:

- 2022-12-09: XML Release v1.0 (TEI encoding) (2 ZIP files)
- 2022-11-01: TXT Release v1.2: remove space before closing bracket.
- 2022-10-07: TXT Release v1.1: correct cite for source text; remove space before end punctuation.
- 2022-10-05: TXT Release v1.01: 3 entries fixed, 2 removed.
- 2022-10-04: We're excited to publish the text of this edition as TXT files.

- **Chapbooks printed in Scotland:** This dataset includes over **3,000 chapbooks** printed in Scotland. Popular reading materials from the late **17th to the 19<sup>th</sup> century**. They were usually printed on a single sheet and then **folded into books of 8, 12, 16 and 24 pages**. **47,329 ALTO XMLs** at page level, ~10 million OCRed words.
- **Ladies' Edinburgh Debating Society :** This collection comprises two Edinburgh journals: *The Attempt* and *The Ladies Edinburgh Magazine*. It spans **16 volumes** from **19<sup>th</sup> century**. **6,354 ALTO XMLs** at the page level ~2.5 million words.
- **Gazetteers of Scotland collection:** **20 volumes** historical gazetteers of Scotland from the **19th century**. **~13,000 ALTO-XMLs** at page level, ~14.5 million words.

New *defoe* heuristic approach – [write\\_metadata\\_pages.yml](#) :

- parallel extraction of text from each page using ALTO-XML files
- extracting metadata from both the collection and volume using METS-XML files

Structuring the extracted information.

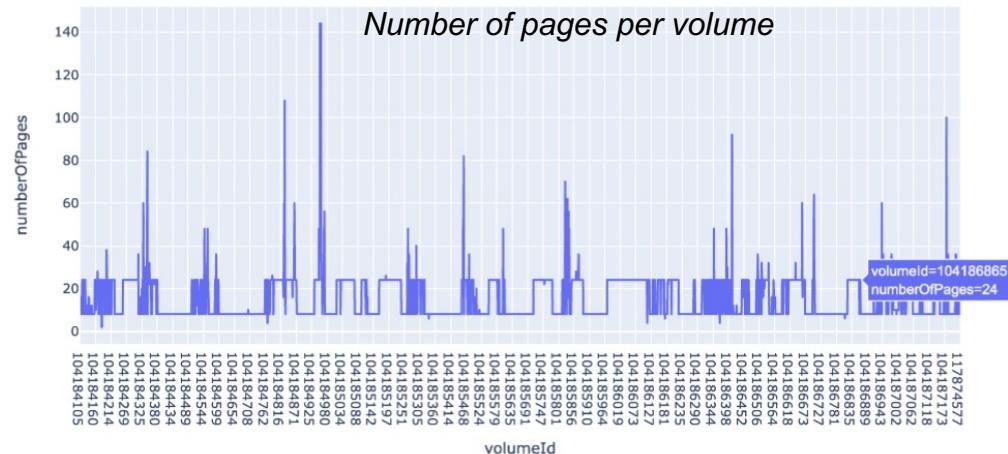
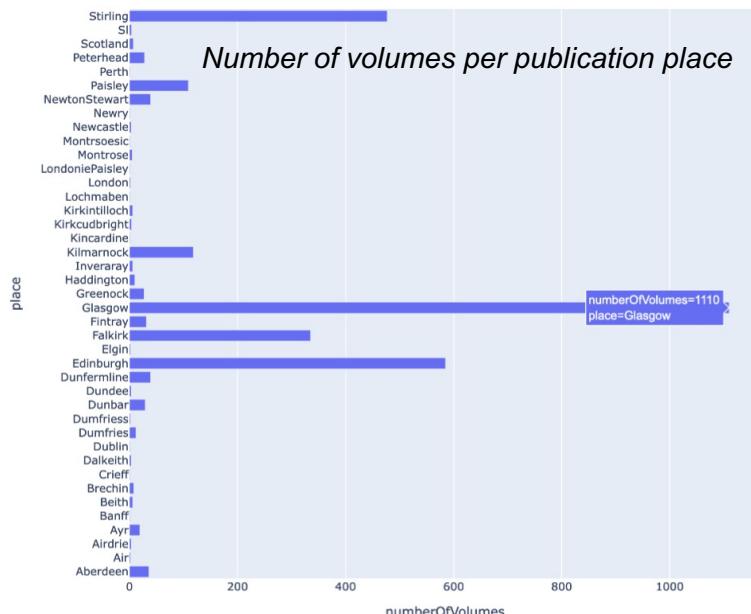
- For each volume - permanent URL that allows visualization of the page images.

# Phase 1.2: Information Extraction – Other NLS Collections

```
MMSID
edition
editor
editor_date
genre
language
metsXML
termsOfAddress
numberOfPages
numberOfWords
permanentURL
physicalDescription
place
publisher
referencedBy
shelfLocator
altoXML
serieSubTitle
text
pageNum
volumeTitle
volumeId
year
serieNum
part
collectionName
serieTitle
publisherPersons
numberOfVolumes
volumeNum
Name: 0, dtype: object
```

**dataframe**

```
9937033633804341
None
Milne, John
1792-1871
Chapbooks-Scotland-Aberdeen-1801-1900
eng
104184105-mets.xml
None
8
53
https://digital.nls.uk/104184105
8 p. ; 18 cm.
Aberdeen
Printed by A. Imlay, 22, Long Acre
None
L.C.2786.A(1)
104184105/alto/107134030.34.xml
to the tune of Johnny Cop
A SONG JRAISB OP THE ^ HIGHLAND LADS. To the T...
1
song in praise of the highland lads
104184105
1826
0
0
Chapbooks printed in Scotland
song in praise of the highland lads
[]
1
1
1
```

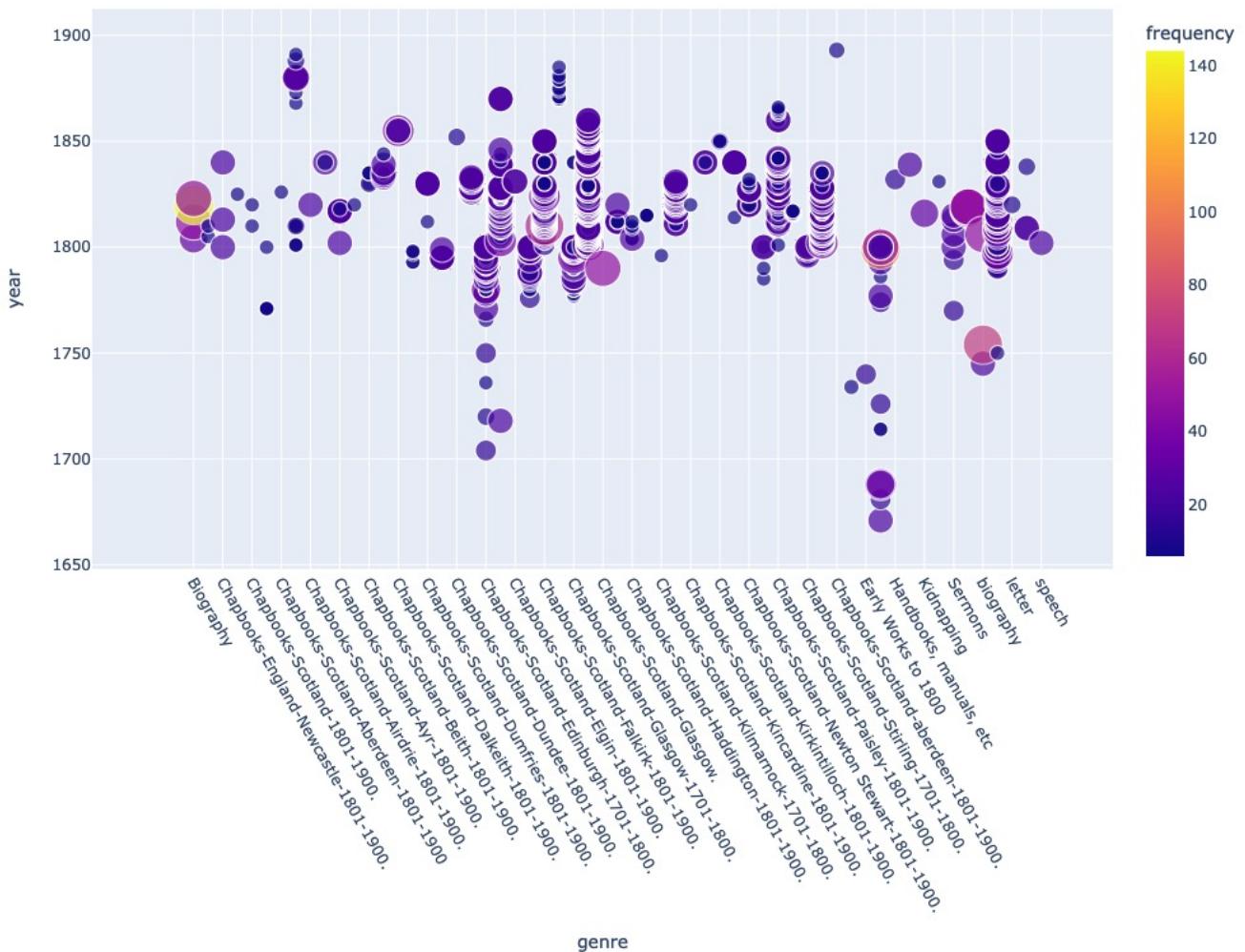


## Phase 1.2: Information Extraction – Other NLS Collections



St Andrews / CS

### Chapbooks frequency per genre and year



*Frequency of genres per year. Each bubble indicates the frequency of a genre for a particular year*

## The Journey

Phase 1: Information Extraction

Phase 2: Ontologies and Knowledge Graphs:

2.1 - EB

2.2 - NLS

Phase 3: Augmented Knowledge Graphs with Deep Transfer Learning

Phase 4: Defoe and Knowledge Graphs

Phase 5: React-Flask Web Platform

Phase 6: Case Studies

## EB Ontology

Aim: to capture a shareable and reusable knowledge representation of the EB.

EB Ontology → A formal description of knowledge as a set of concepts

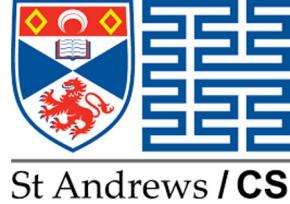
In order to create and publish the EB-Ontology we used:

- diagrams.net : To create an UML with the EB information (classes, properties, relationships, etc.)
- Chowlk : To convert the UML into an OWL ontology
- Widoco: To publish and create an enriched and customized documentation of the ontology
- w3id.org: To configure my permanent Identifier for EB ontology

EB-Ontology available at:

- <https://github.com/francesNLP/EB-ontology>
- <https://w3id.org/eb/>

# Phase 2.1: EB Ontology and Knowledge Graph



## EB Ontology

This version:

<https://w3id.org/eb/0.0.2>

Latest version:

<https://w3id.org/eb>

Revision:

0.0.2

Authors:

Rosa Filgueira

Contributors:

Daniel Garijo

Download serialization:

[Format](#) [JSON LD](#) [Format](#) [RDF/XML](#) [Format](#) [N Triples](#) [Format](#) [TTL](#)

License:

[License](#) license name goes here

Visualization:

[Visualize with](#) [WebVowl](#)

Cite as:

Rosa Filgueira. EB Ontology. Revision: 0.0.2. Retrieved from: <https://w3id.org/eb/0.0.2>

## Abstract

An ontology designed to represent volumes, terms and editions of the Encyclopaedia Britannica. The ontology extends Schema.org.

## Table of contents

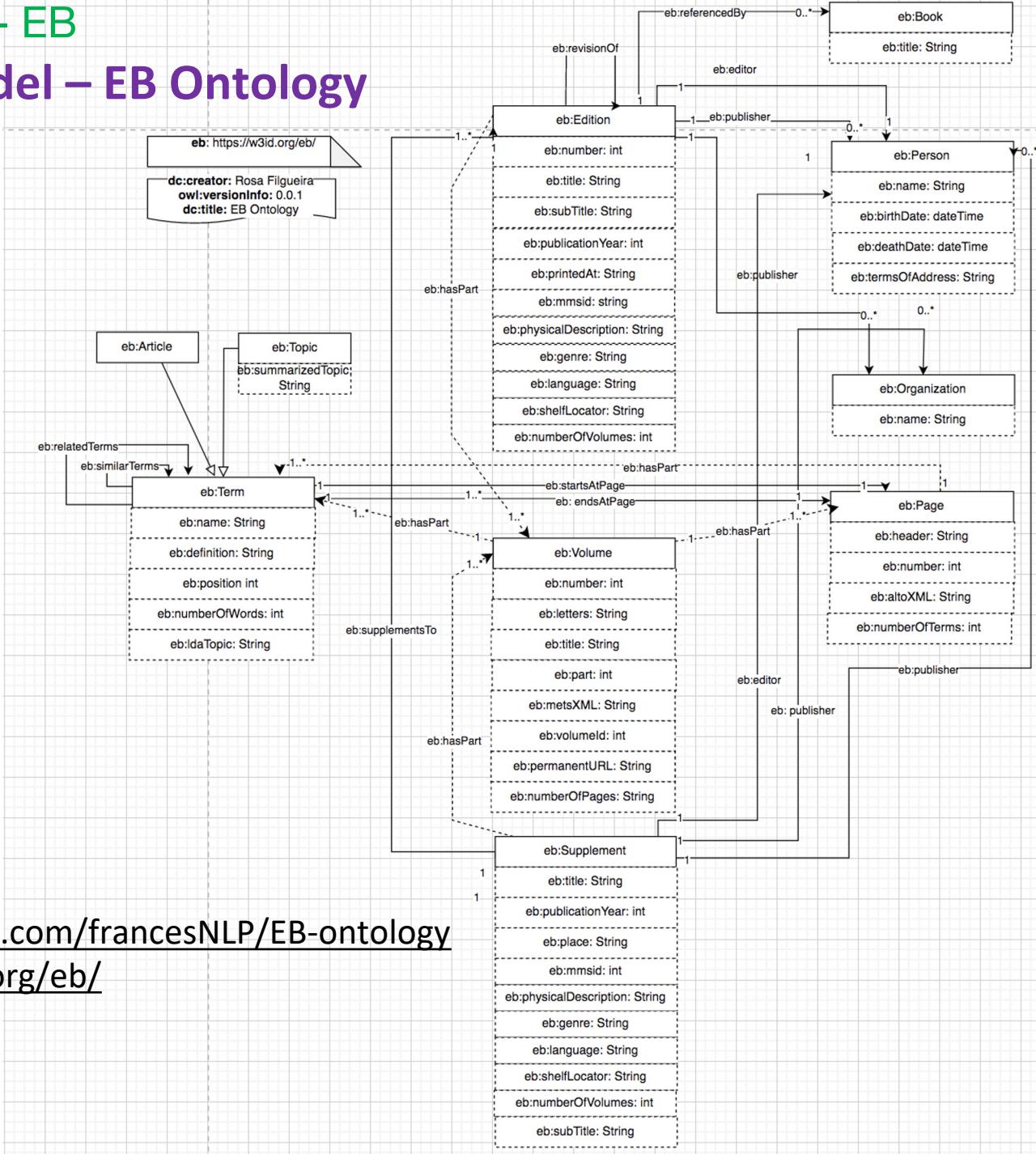
- 1. [EB Ontology: Description](#)
- 2. [Cross reference for EB Ontology classes, properties and dataproperties](#)
  - 2.1. [Classes](#)
  - 2.2. [Object Properties](#)
  - 2.3. [Data Properties](#)
- 3. [Acknowledgments](#)

## 1. Introduction

This is a place holder text for the introduction. The introduction should briefly describe the ontology, its motivation, state of the art and goals.

# Phase 2.1 - EB

## Data Model – EB Ontology



<https://github.com/francesNLP/EB-ontology>

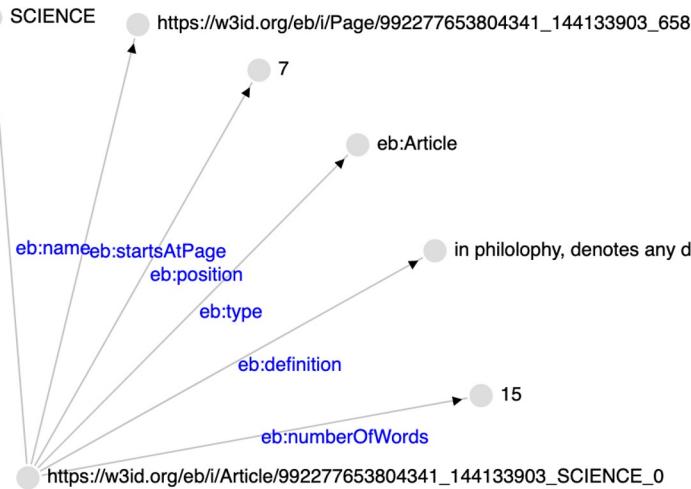
<https://w3id.org/eb/>

### EB Knowledge Graph

**EB Knowledge Graph:** Populating the post-processed information (extracted **Terms & Metadata**) into an RDF triple store + **EB Ontology**:

- Apache Jena FUSEKI SPARQL server

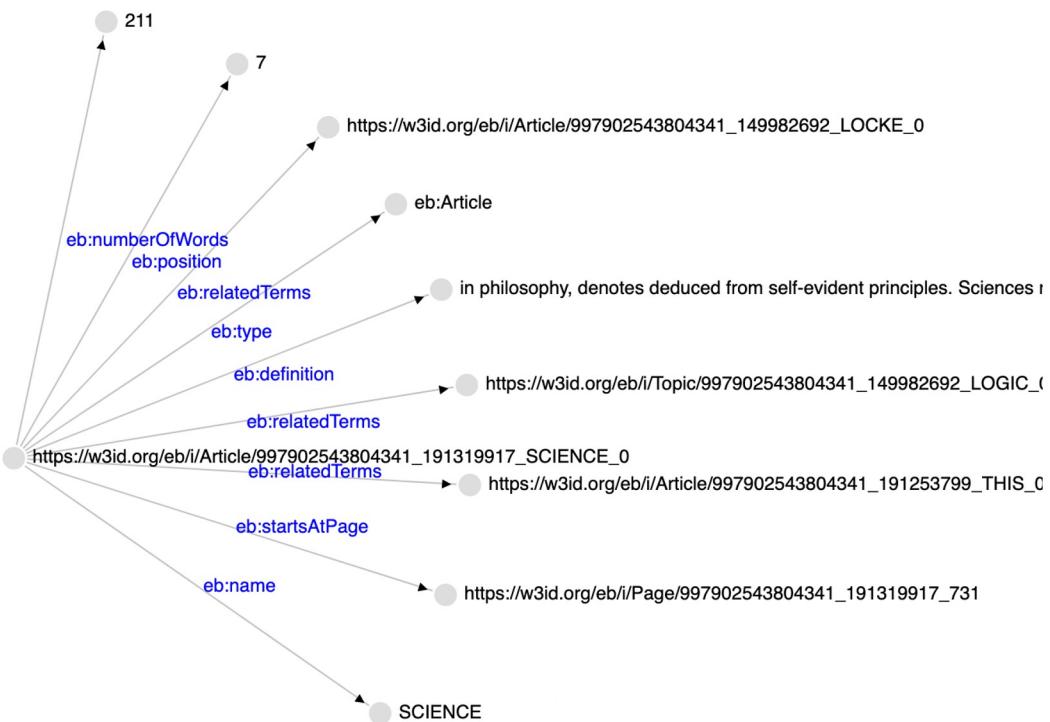
**NOTE:** Each Term, Edition, Page, Volume, etc ... is a **Resource in our Knowledge Graph** and has an **URI to identify it**.



- **Science** term in the First Edition issued in 1771

Edition 1 Year 1771  
 Edition 1 Year 1773  
 Edition 2 Year 1778  
 Edition 3 Year 1797  
 Edition 4 Year 1810  
 Edition 5 Year 1815  
 Edition 6 Year 1823  
 Edition 7 Year 1842  
 Edition 8 Year 1853

---



- **Science** term in the Third Edition issued in 1797

# Phase 2.1: EB Ontology and Knowledge Graph



zenodo

Search

Upload Communities

There is a **newer version** of this record available.

June 21, 2022

Dataset Open Access

## EB-KG: Knowlege Graph of the first 8 eiditions Encyclopaedia Britannica (1768-1860)

Rosa Filgueira

This Knowlege Graph represents the information of the first eight editions of Encyclopaedia Britannica (years: 1768 to 1860) in RDF (ttl format).

The raw dataset is provided by the NLS in this [link](#), and it comprises of eight editions and a total of 195 volumes with a total size of 44GB. It uses two XMLS schemas: METS for descriptive, structural, technical and administrative metadata (Title, Author, Publisher, etc); and ALTO for encoding the OCR text of a page.

In this work, we have extracted the information from METS and ALTO XMLS using [defoe](#) tool and developed [novel information extraction heuristics](#). With the extracted information, we created the EB-KG Knowlege Graph, which uses the [EB Ontogy](#), to represent such information. Furthermore, during the information extraction phase, we have employed several techniques to mitigate two common OCR errors: long-S and the line-break hyphenation.

The EB-KG contains 1,638,239 RDF triples. It has information from 8 editions. Each edition can have several Volumes, references to Books, Supplements; it also has an Editor and a Publisher, which can be a Person or an Organization. A Volume has several Pages, which can contain several Terms. And a Term can be either a Topic (a term described across several pages, often combining text, pictures, and tables.) or an Article (a description of the term in one- or two-paragraph long text (similar to an entry in a dictionary)). The data model of the EB-KG can be found [here](#).

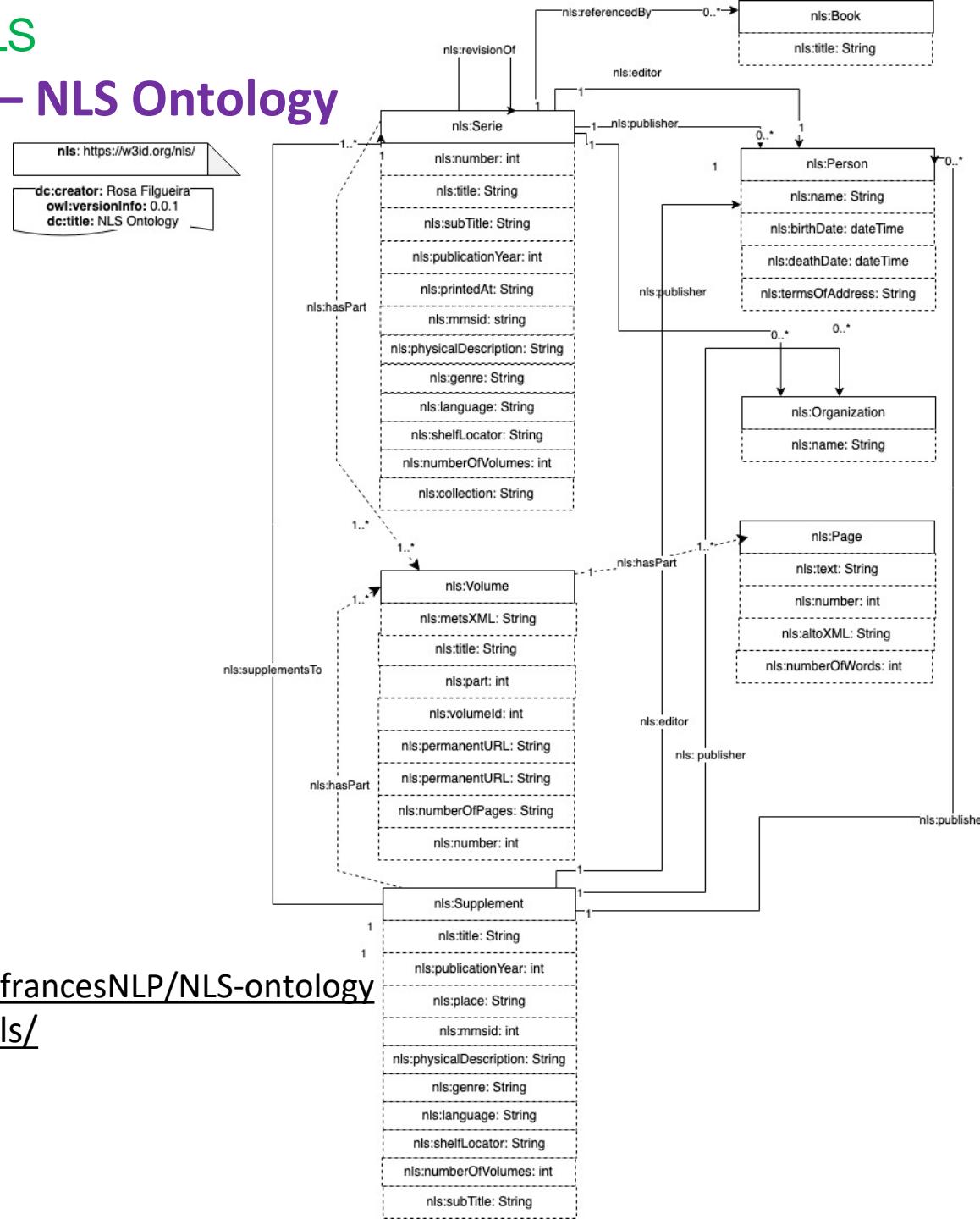
The original ALTO files do not indicate the start and end of each EB term, the first part of our work involved the automated extraction of all terms (along with their metadata) across editions, so they can be analysed independently without the surrounding text.

This work was performed during my 2021-2022 National Library of Scotland Digital Scholarship Fellowship.

## Phase 2.2 - NLS



# Data Model – NLS Ontology



<https://github.com/francesNLP/NLS-ontology>

<https://w3id.org/nls/>

### EB Knowledge Graph

**3 NLS Knowledge Graph:** Populating the post-processed information into RDF triple store + **NLS Ontology:** Apache Jena FUSEKI SPARQL server

[gazetters\\_scotland.ttl](#)

[chapbooks\\_scotland.ttl](#)

[ladies\\_debating.ttl](#)

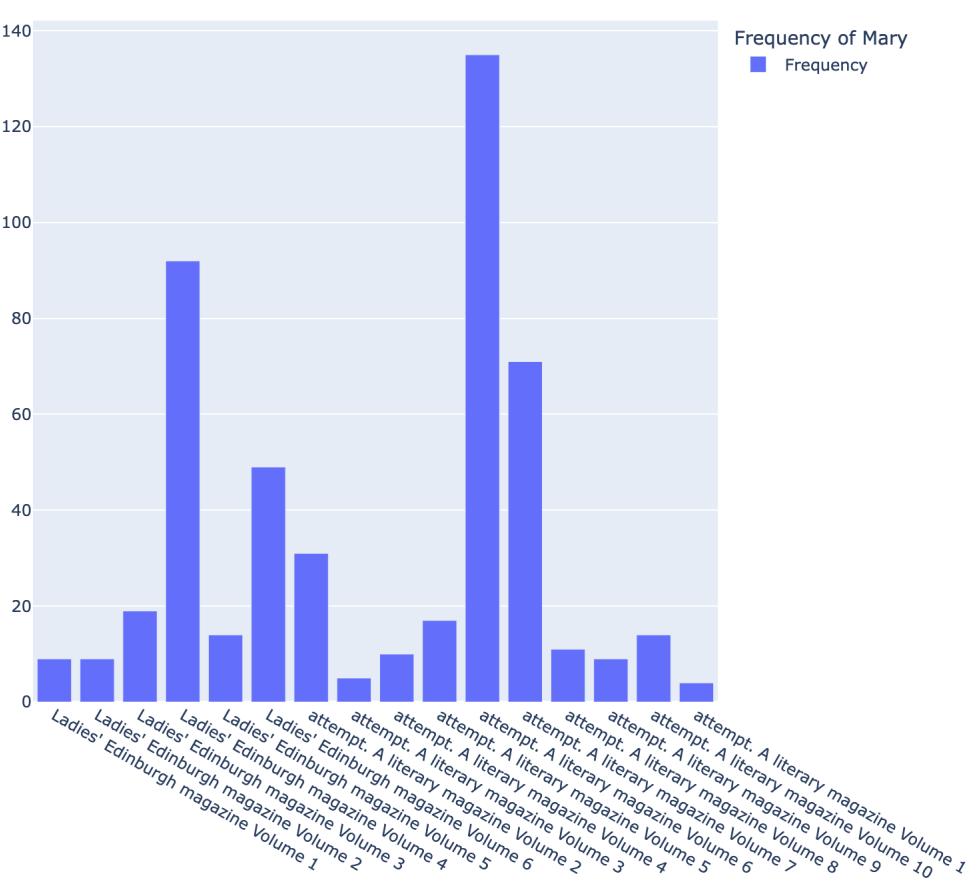
## Phase 2.2: NLS Ontology and Knowledge Graphs

Calculating the frequency of the term ‘Mary’ across the publications of the *Ladies’ Edinburgh Debating Society* using [ladies\\_debating.ttl](#)

```
sparql = SPARQLWrapper("http://35.228.63.82:3030/ladies_debating/sparql")
sparql.setQuery("""
SELECT ?text ?title
WHERE {
?page a nls:Page .
?page nls:text ?text .
?v nls:hasPart ?page .
?v nls:title ?title
FILTER regex(?text, "Mary")
}

""")
sparql.setReturnFormat(JSON)
results = sparql.query().convert()
```

SPARQL Query to obtain the frequency of ‘Mary’ across volumes.



Visualization of the frequency results calculated above.

# Phase 2.2: NLS Ontology and Knowledge Graphs



zenodo

Search

Upload Communities

Log in Sign up

June 22, 2022

## GazetteersScotland-KG: A Knowledge Graph for representing the Gazetteers of Scotland (1803-1901)

Rosa Filgueira

This Knowledge Graph represents the information of the "Gazetteers of Scotland" (years: 1803 - 1901) collection in RDF (ttl format). This collection comprises twenty volumes of the most popular descriptive gazetteers of Scotland in the 19th century. Principal places in Scotland, including towns, counties, castles, glens, antiquities and parishes, are listed alphabetically. Each entry includes detailed historical and geographical information about each place. The raw dataset is provided by the NLS in this [link](#). As other NLS data collections, they are originally provided using two XMLS schemas: METS for descriptive, structural, technical and administrative metadata (Title, Author, Publisher, etc); and ALTO for encoding the OCR text of a page.

In this work, we have extracted the information from METS and ALTO XMLS using [defoe](#) tool and developed a new [information extraction defoe query](#), and created a new Knowledge Graph called GazetteersScotland-KG. The GazetteersScotland-KG uses the [NLS Ontology](#) to represent the information extracted. Furthermore, during the information extraction phase, we have employed several techniques to mitigate two common OCR errors: long-S and the line-break hyphenation.

The GazetteersScotland-KG contains 354,998 RDF triples. It has information from 12 series and 20 volumes: Each serie can have several Volumes. Each serie has an Editorm Publisher, mmsid, Shelf-Locator, publication year, etc. A Volume has several Pages, with text in them. The data model of the GazetteersScotland-KG can be found [here](#).

Dataset Open Access

zenodo

Search

Upload Communities

June 22, 2022

## LadiesDebating-KG: A Knowledge Graph for representing the "Edinburgh Ladies' Debating Society Digital Collection" (1865 - 1880)

Rosa Filgueira

This Knowledge Graph represents the information of the "Edinburgh Ladies' Debating Society" (years: 1865 - 1880) collection in RDF (ttl format). This collection consists of the complete runs of two Edinburgh journals, **'The Attempt'** (10 volumes, 1865-74) and its successor **'The Ladies' Edinburgh Magazine'** (6 volumes, 1875-80). These publications were produced by a leading Edinburgh women's club, known during the period as the Edinburgh Essay Society or the Ladies' Edinburgh Essay Society, but subsequently as the Ladies' Edinburgh Debating Society. The Society existed from 1865 to 1935. The raw dataset is provided by the NLS in this [link](#). As other NLS data collections, they are originally provided using two XMLS schemas: METS for descriptive, structural, technical and administrative metadata (Title, Author, Publisher, etc); and ALTO for encoding the OCR text of a page.

In this work, we have extracted the information from METS and ALTO XMLS using [defoe](#) tool and developed a new [information extraction defoe query](#), and created a new Knowledge Graph called LadiesDebating-KG. The LadiesDebating-KG uses the [NLS Ontology](#) to represent the information extracted. Furthermore, during the information extraction phase, we have employed several techniques to mitigate two common OCR errors: long-S and the line-break hyphenation.

The LadiesDebating-KG contains 38,279 RDF triples. It has information from 2 series and 16 volumes: **'The attempt'** serie has 10 volumes and **'The Ladies'** serie has 6 volumes. Each serie has an Editor, mmsid, Shelf-Locator, publication year, etc. A Volume has several Pages, with text in them. The data model of the LadiesDebating-KG can be found [here](#).

Dataset Open Access

zenodo

Search

Upload Communities

Log in Sign up

30 views 30 views

Indexed in  
OpenAIRE

Publication date: June 22, 2022  
DOI: DOI 10.5281/[zenodo.6696686](#)  
Keyword(s): Knowledge Graph, Semantic Web, RDF, Chapbooks Printed In Scotland, Digital Humanities  
License (for files): CC-BY Creative Commons Attribution 4.0 International

Files (75.5 MB)  
Name Size  
chapbooks\_scotland.ttl 75.5 MB  
[Download](#)

43 views 2 downloads  
See more details...

Indexed in  
OpenAIRE

Publication date: June 22, 2022  
DOI: DOI 10.5281/[zenodo.6686596](#)  
Keyword(s): Knowledge Graph, Semantic Web, RDF, Chapbooks Printed In Scotland, Digital Humanities  
License (for files): CC-BY Creative Commons Attribution 4.0 International

96 views 6 downloads  
See more details...

Indexed in  
OpenAIRE

Publication date: June 21, 2022  
DOI: DOI 10.5281/[zenodo.6696686](#)  
Keyword(s): Knowledge Graph, Semantic Web, RDF, Chapbooks Printed In Scotland, Digital Humanities  
License (for files): CC-BY Creative Commons Attribution 4.0 International

Versions

## The Journey

Phase 1: Information Extraction

Phase 2: Ontologies and Knowledge Graphs

**Phase 3: Augmented Knowledge Graph with Deep Transfer Learning**

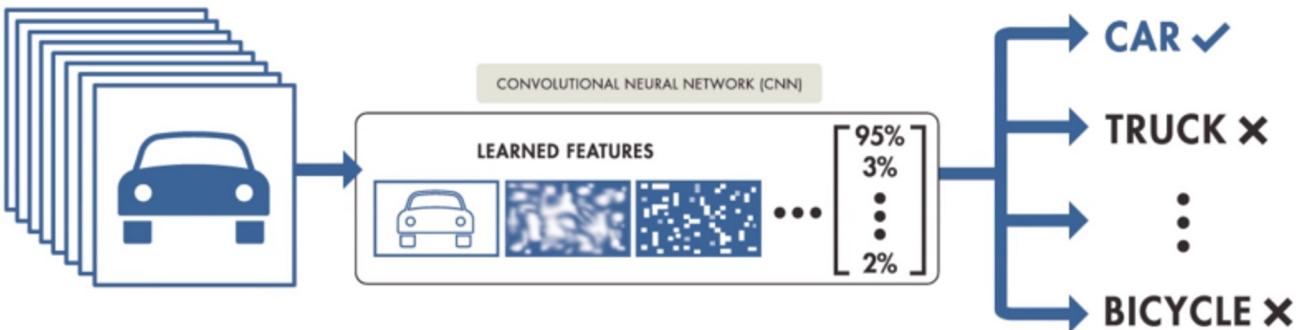
Phase 4: Defoe and Knowledge Graphs

Phase 5: React-Flask Web Platform

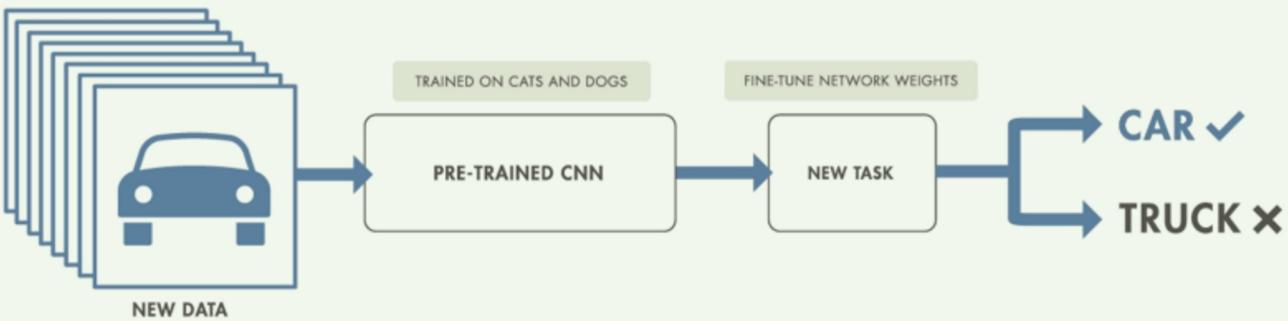
Phase 6: Case Studies

## Background

### TRAINING FROM SCRATCH



### TRANSFER LEARNING



We are going to use this approach for applying NLP/ML analysis

## Background



# Transformers

build passing

license Apache-2.0

website online

release v2.0.0

## State-of-the-art Natural Language Processing for PyTorch and TensorFlow 2.0

😊 Transformers provides thousands of pretrained models to perform tasks on texts such as classification, information extraction, question answering, summarization, translation, text generation, etc in 100+ languages. Its aim is to make cutting-edge NLP easier to use for everyone.

😊 Transformers provides APIs to quickly download and use those pretrained models on a given text, fine-tune them on your own datasets then share them with the community on our [model hub](#). At the same time, each python module defining an architecture can be used as a standalone and modified to enable quick research experiments.

😊 Transformers is backed by the two most popular deep learning libraries, [PyTorch](#) and [TensorFlow](#), with a seamless integration between them, allowing you to train your models with one then load it for inference with the other.

## Background



build passing license Apache-2.0 website online release v2.0.0

### A High-Level Look

Let's begin by looking at the model as a single black box. In a machine translation application, it would take a sentence in one language, and output its translation in another.



## Augmenting the EB-KG

**EB Knowledge Graph 2.0:** previous info + storing the result of applying different deep learning transformers analyses:

- **sentiment analyses:** Classifying articles between positive and negative
  - **transformer:** siebert/sentiment-roberta-large-english
- **topic modelling:** Clustering terms into topics
  - **transformer:** all-mpnet-base-v2
- **term similarity:** Comparing terms & semantic similarity
  - **transformer:** all-mpnet-base-v2
- **spelling checking:** Finding misspelling/ocr errors and fixing them
  - **transformer:** neuspell + ElmoslstmChecker
- **summarization:** Summarizing the text of a topic term (XLNET)
  - **transformer:** XLNet

## Augmenting the EB-KG

**Example:** Spelling Checking --> *Lewis* Term

### Original Definition

the most northerly of any of the western islands of Scotland, lying in 8° odd minutes W. long, and between 58° and 59 0 odd minutes N. lat.



### Cleaned Definition

the most northerly of any of the western islands of Scotland , lying in 8 and odd minutes W. long , and between 58 and and 59 0 odd minutes N. land .



### Compute Difference

the morst northerly of any of the w easter islands of Scotland, lying in 8\u00b0 and odd minutes W. long, and between 58\u00b0 and 59 0 odd minutes N. latnd.

## The Journey

Phase 1: Information Extraction

Phase 2: Ontologies and Knowledge Graphs

Phase 3: Augmented Knowledge Graphs with Deep Transfer Learning

**Phase 4: Defoe and Knowledge Graphs**

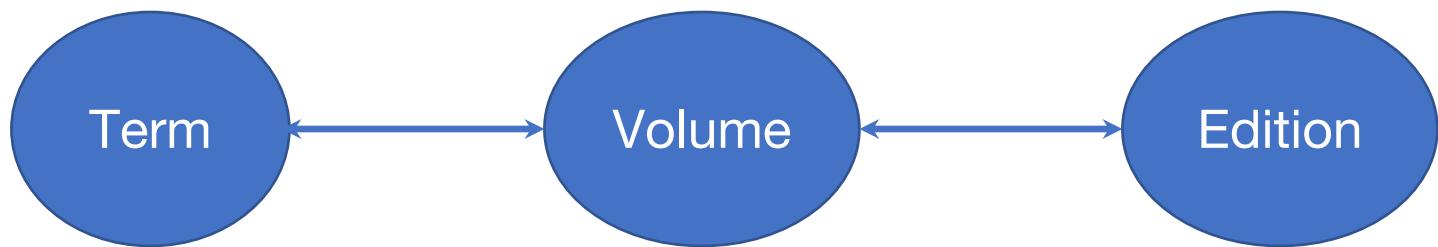
Phase 5: React-Flask Web Platform

Phase 6: Case Studies

## Querying the KGs

Two types of “queries”:

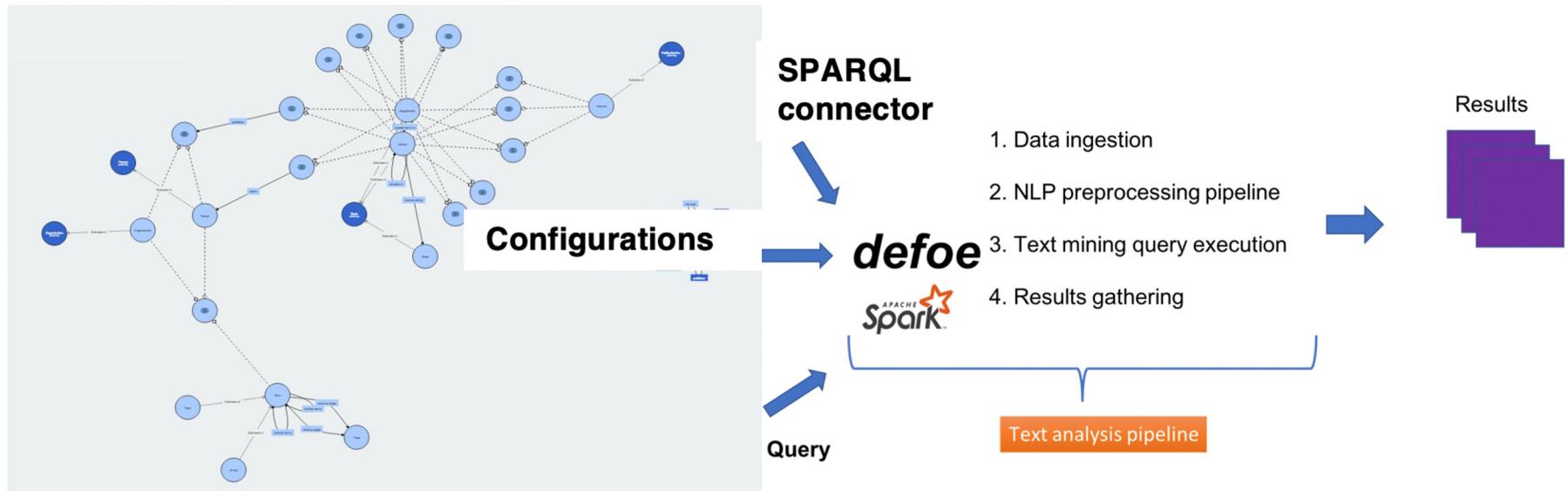
- **Type 1: Extracting information from the EB Knowledge → SPARQL:**  
Navigating through the KG to get the desired information.



Example: Given a **Term** ( e.g., *Edinburgh*), I can get the **Edition** Information (e.g., *Edition Title*)

- **Type 2: Processing information from the Knowledge Graphs → *defoe*:**  
Processing further the definitions from the selected terms/pages.
  - We needed to do some work on *defoe* first --> **Phase 4**

New **KG defoe connector** (based in SPARQL) to run *defoe* queries using the EB Knowledge Graph as a source of data → **defoe\_KG**



New *defoe* queries fully configurable: different filtering options, target, lexicon, etc.

- frequency keysearch: **Counts number of terms/pages or times** in which appear keywords or keysentences and group by years. Several filtering options.
- fulltext keysearch: **Extracts terms definitions/pages** in which appear keywords or keysentences and group by years. Several filtering options.
- snippet keysearch: **Extracts snippets of definitions/pages** in which appear keywords or keysentences and group by years. Filtering options, including the snippet size.
- publication normalization: **Extracts the number of documents, pages, words** per year.
- uris keysearch: **Extract uris of terms/pages** in which appear keywords or keysentences and group by years. Several filtering options.
- geoparser terms: **Geoparsing the term definition/pages** in which appear keywords or keysentences and group by years. Several filtering options.
- frequency-distribution: Calculates the **frequency of the most 'N' common tokens** in terms definitions/pages.
- lexicon-diversity: Computes the **lexical diversity** metric for a given collection: ratio of the vocabulary size to the total number of words in the text.
- person entity recognition: **Identifies people mentioned** in text and estimates the gender distribution.

## The Journey

Phase 1: Information Extraction

Phase 2: Ontologies and Knowledge Graphs

Phase 3: Augmented Knowledge Graphs with Deep Transfer Learning

Phase 4: Defoe and Knowledge Graphs

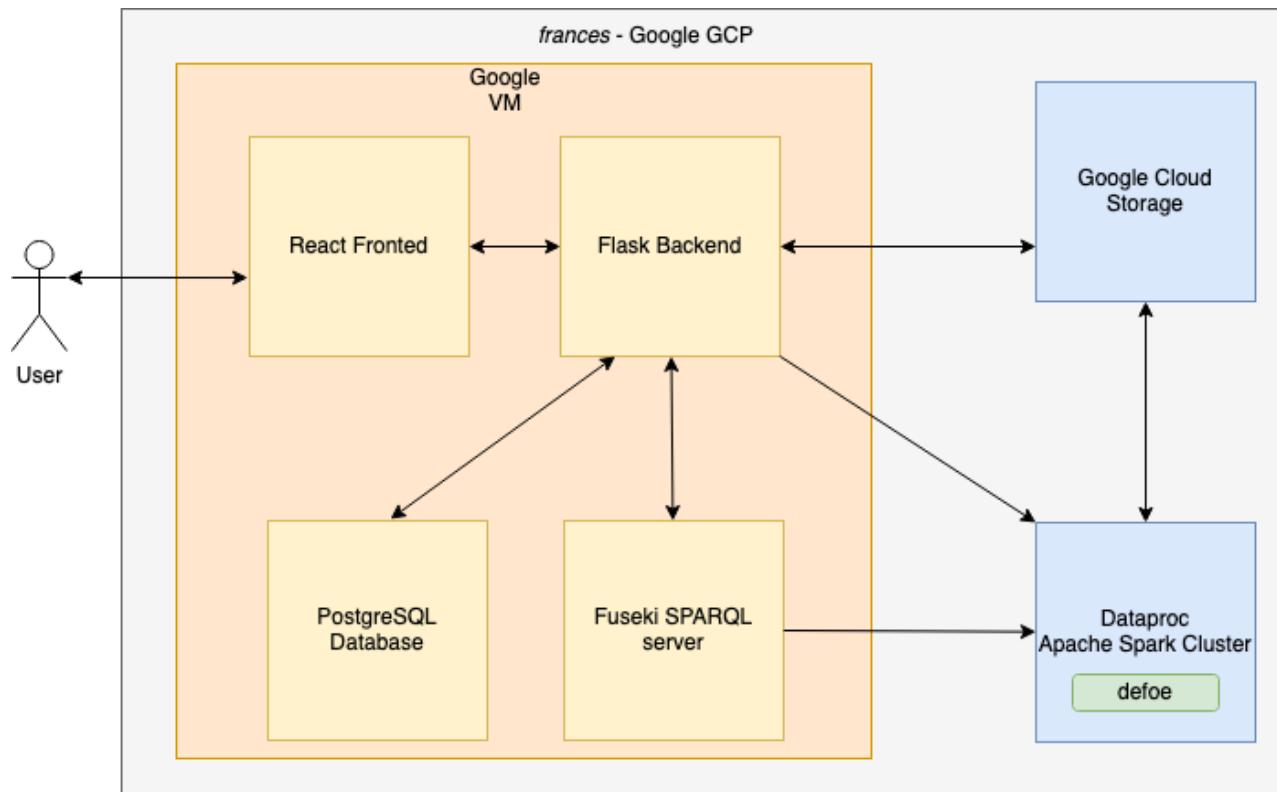
**Phase 5: React-Flask Web Platform**

Phase 6: Case Studies

## Phase 5: React-Flask Web Platform



- Users DO NOT have to create SPARQL or *defoe* queries → the web-application does it for them → Abstractions to SPARQL & DEFOE
- The web-app runs both types of queries and visualizes the results



← → ⌂

www.frances-ai.com

☆ ⌂ ⌄ ⌅ ⌆ ⌇ ⌈ ⌉

frances Term Search Term Similarity Topic Modelling Defoe Queries Collection Detail CREATE AN ACCOUNT SIGN IN

## Exploring the Encyclopaedia Britannica (1768-1860)

### Term Search

Enter the **term** that you would like to search for. In case that the **Term Type is a Topic**, only the **summary** of the definition is displayed. If not term is introduced, it will search for the first term in the Encyclopaedia.

# Phase 5: React-Flask Web Platform



frances

Term Search

Term Similarity

Topic Modelling

Defoe Queries

Collection Details

CREATE AN ACCOUNT

SIGN IN

Search Term

Flower



Results for FLOWER

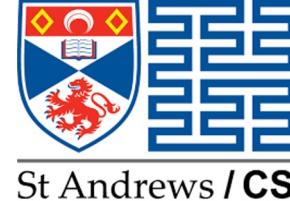
Note that if you click over a **related term**, it will conduct a **term search**, showing all the searching results for that term. And if you click over a **topic model id** if will take you to the **Topic Modelling page**, listing all the terms belonging to that particular topic model.

Year	Edition	Volume	Start Page	End Page	Term Type	Definition / Summary	Related Terms	Topic Modelling ID	Sentiment Score	Advanced Options
1771	1	2	566	566	Article	among botaniffs and gardeners, the most beautiful part of trees and plants, containing the organs or parts of frutification. See Botany., Fr/cr W Flower. See Xeranthemum. Everlajiinx Flower.. See Gna... ↗		3190	POSITIVE_1.00	<a href="#">VISUALISE</a> <a href="#">CHECK</a> <a href="#">SPELL</a> <a href="#">SIMILAR TERMS</a>
1773	1	2	574	574	Article	among botaniffs and gardeners, the most beautiful part of trees and plants, containing the organs or parts of frutification. See Botany., External Flower. See Xeranthemum, Everlafling Flower. See Gna... ↗	BOTANY	3190	POSITIVE_1.00	<a href="#">VISUALISE</a> <a href="#">CHECK</a> <a href="#">SPELL</a> <a href="#">SIMILAR TERMS</a>

1–2 of 9



# Phase 5: React-Flask Web Platform



frances

Term Search

Term Similarity

Topic Modelling

Defoe Queries

Collection Details

CREATE AN ACCOUNT

SIGN IN

## Term Similarity

Enter **some text** that you would like to search similar terms for. If no term is introduced, it will search for the first term in the Encyclopaedia.


The first **20** most similar results. The results are sorted by **similitud rank**.

Note that if you click over a **related term**, it will conduct a **term search**, showing all the searching results for that term. And if you click over a **topic model id** if will take you to the **Topic Modelling page**, listing all the terms belonging to that particular topic model.

Year	Edition	Volume	Term	Definition	Topic Modelling ID	Similitud rank	Sentiment Score	Advanced Options
1823	6	8	FLOWER	Flos, among botanists and gardeners, the most beautiful part of trees and plants, containing the organs or parts of fructification. Seedesigned for medical use, should be jilnck— ■ / gj when they are ... ↗	-1	0.5790331363677979	LABEL_0_0.87	<a href="#">VISUALISE</a> <a href="#">CHECK</a> <a href="#">SPELL</a> <a href="#">SIMILAR TERMS</a>
1823	6	16	PLENUM	FLOS, a full flower 5 a the highest degree of luxuriance in flowers. tany. Such flowers, although the the eye, are both vegetable monsters, to the sexualists, vegetable eunuchs; crease of the petals c... ↗	1231	0.5691936016082764	LABEL_0_0.91	<a href="#">VISUALISE</a> <a href="#">CHECK</a> <a href="#">SPELL</a> <a href="#">SIMILAR TERMS</a>

# Phase 5: React-Flask Web Platform



frances

Term Search

Term Similarity

Topic Modelling

Defoe Queries

Collection Details

CREATE AN ACCOUNT

SIGN IN

Place some text

person who does scientific experiments

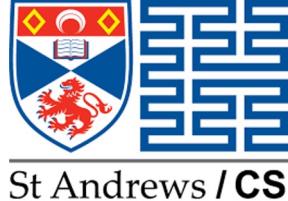


The first **20** most similar results. The results are sorted by **similitud rank**.

Note that if you click over a **related term**, it will conduct a **term search**, showing all the searching results for that term. And if you click over a **topic model id** it will take you to the **Topic Modelling page**, listing all the terms belonging to that particular topic model.

Year	Edition	Volume	Term	Definition	Topic Modelling ID	Similitud rank	Sentiment Score	Advanced Options
1823	6	2	ARMORIST	a person skilled in the knowledge of	<a href="#">3715</a>	0.5842932462692261	POSITIVE_1.00	<a href="#">VISUALISE</a> <a href="#">CHECK</a> <a href="#">SPELL</a> <a href="#">SIMILAR TERMS</a>
1842	7	18	PROFESSOR	in the universities, a person who teaches or reads public lectures in some art or science. See Uni	<a href="#">441</a>	0.5636317133903503	POSITIVE_0.99	<a href="#">VISUALISE</a> <a href="#">CHECK</a> <a href="#">SPELL</a> <a href="#">SIMILAR TERMS</a>
1815	5	16	PHYSICIAN	a person who professes medicine, or the art of healing diseases. See Medicine. PII YSI CI AKS, College of, in London, Edinburgh, Dublin.	<a href="#">3251</a>	0.5457887649536133	POSITIVE_1.00	<a href="#">VISUALISE</a> <a href="#">CHECK</a> <a href="#">SPELL</a>

# Phase 5: React-Flask Web Platform



## Topic Modelling Search

Enter a **topic name** or just the **ID** of a topic model to see all the terms within the same topic. All topics modelling names start with a number. If no topic is introduced, it will use the first topic modelling, which name starts with '0'.

11 terms found for the topic **3715\_knowledge\_acquired\_any\_teacher\_use\_without\_teacher\_mr**

Note that if you instead click over a **term**, it will take you to the [Term Search page](#), showing all the searching results for that term.

Year	Edition	Volume	Term	Definition	Sentiment Score	Advanced Options
1810	4	16	<a href="#">PHIOMATHES</a>	a lover of learning or fciaeace. thes,	POSITIVE_1.00	<a href="#">VISUALISE</a> <a href="#">CHECK SPELL</a> <a href="#">SIMILAR TERMS</a>
1823	6	2	<a href="#">ARMORIST</a>	a person skilled in the knowledge of	POSITIVE_1.00	<a href="#">VISUALISE</a> <a href="#">CHECK SPELL</a> <a href="#">SIMILAR TERMS</a>
1815	5	2	<a href="#">ARMORIST</a>	a person killed in the knowledge of	NEGATIVE_1.00	<a href="#">VISUALISE</a> <a href="#">CHECK SPELL</a> <a href="#">SIMILAR TERMS</a>

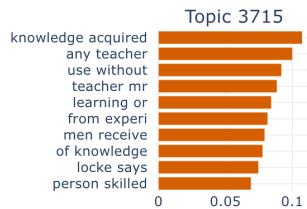
## Topic Modelling Visualisations

We are visualising each topic model found in our term-search. For each topic model, we plot their most common 10 words , along with their [c-TF-IDF scores](#).

Note that the topic that starts with "-1" refers to all outliers and should typically be ignored



### Topic Word Scores



# Phase 5: React-Flask Web Platform

frances

Term Search

Term Similarity

Topic Modelling

Defoe Queries

Collection Details

CREATE AN ACCOUNT

SIGN IN

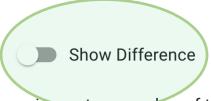
TermSearch(ABACUS) > SpellCheck(992277653804341\_144133901\_ABACUS\_1)

Year	Edition	Volume	Term	Advanced Options
1771	1	1	ABACUS	<a href="#">VISUALISE</a> <a href="#">SIMILAR TERMS</a>

## Original Definition

in architecture, signifies the superior part or member of the capital of a column, and serves as a kind of crowning to both. It was originally intended to represent a square tile covering a basket. The form of the abacus is not the same in all orders: in the Tuscan, Doric, and Ionic, it2018is generally square; but in the Corinthian and Compoite, its four sides are arched ir Avars, and embellilhed in the middle withornament, as a rose or other flower, Scammozzi uses abacus for a concave moulding on the capital of the Tuscan pedefial; and Palladio calls the plinth above the echinus, or boultin, in the Tufean and Doric orders, by the same name. See plate I. fig. i. and

## Cleaned Definition



in architecture , signifies the superior part or member of the capital of a column , and serves as a kind of crowning to both . It was originally intended to represent a square tile covering a basket . The form of the abacus is not the same in all orders : in the Tuscan , Doric , and Ionic , its generally square ; but in the Corinthian and Compoite , its four sides are arched in Avars , and embellished in the middle withornament , as a rose or other flower , Scammozzi uses abacus for a concave moulding on the capital of the Tuscan pedestal ; and Palladio calls the plinth above the echinus , or boultin , in the Tufean and Doric orders , by the same name . See plate I. fig . i. and



## Cleaned Definition



in architeflure, signifies the superior part or member of the capital of a column, and serves as a kind of crowning to both. It was originally intended to represent a square tile covering a basket. The form of the abacus is not the same in all orders; in the Tuscan, Doric, and Ionic, it2018is generally square; but in the Corinthian and Compoite, its four sides are arched in Avars, and embellished in the middle withornament, as a rose or other flower, Scammozzi uses abacus for a concave moulding on the capital of the Tuscan pedefistal; and Palladio calls the plinth above the echinus, or boultin, in the Tufean and Doric orders, by the same name. See plate I. fig. i. and

frances

Term Search

Term Similarity

Topic Modelling

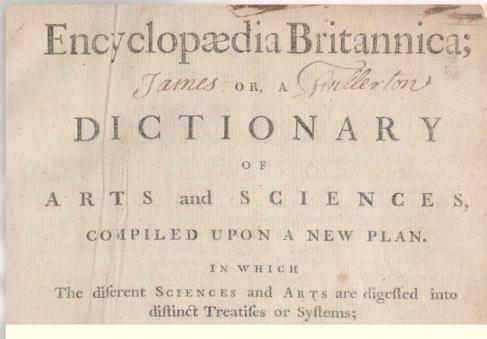
Defoe Queries

Collection Details

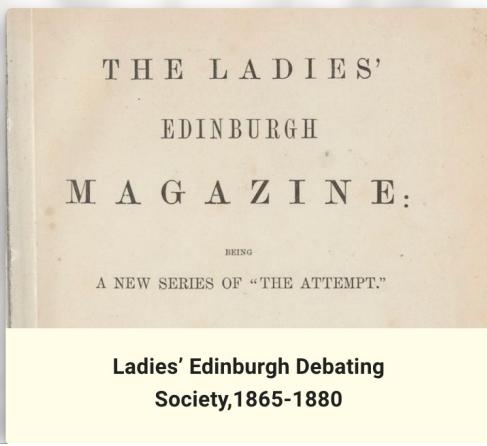
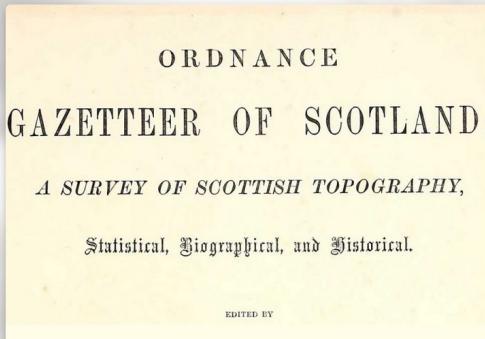
CREATE AN ACCOUNT

SIGN IN

## Exploring NLS Digital Collections



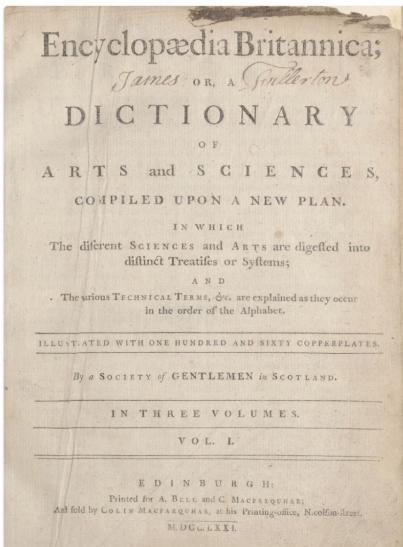
Chapbooks printed in Scotland, 1700-1899



# Phase 5: React-Flask Web Platform



St Andrews / CS



Covers years **1768-1860**

Number of image files **155388**

SOURCE

## Encyclopaedia Britannica

The first eight editions of Encyclopaedia Britannica, issued from 1768-1860, comprise a total of 143 volumes. The Britannica was first issued in Edinburgh in 100 weekly parts (forming 3 volumes) from 1768 to 1771 and illustrated with 160 copperplate engravings. The enterprise was undertaken by the partnership of printer Colin Macfarquhar (1744-1793) and engraver Andrew Bell (1726-1809) who paid William Smellie (1740-1795) to compile the Britannica's first edition for a fee of £200.

Subsequent editions of the Britannica expanded the content: the second edition was published in 10 volumes (1777-1784); the third in 18 volumes (1788-1797). As the Britannica expanded it sought contributions from leading experts in their fields who were either approached by the editors or drawn by the encyclopaedia's growing reputation. Macfarquhar and Bell retained the copyright for the first three editions before publication was taken over by Archibald Constable and then A & C Black. Managed and published in Edinburgh up to the 9th edition.

The Britannica set the standard for modern encyclopedias and is sometimes seen as an enduring product of the Scottish Enlightenment. These volumes were a compendium of current and practical knowledge made relatively affordable by the initial efforts of Macfarquhar and Bell and by Smellie's views on the democratisation of knowledge and the axiom, with which he opens the Preface to the first edition, that "utility ought to be the principal intention of every publication."

### Editions for Encyclopaedia Britannica

Search Edition

- [Edition 1,1771](#)
- [Edition 1,1773](#)
- [Edition 2,1778](#)
- [Edition 3,1797](#)
- [Sup. Edition 3, 1801](#)
- [Edition 4,1810](#)
- [Edition 5,1815](#)

Title: **Edition 1,1771**

SubTitle: Illustrated with one hundred and sixty copperplates

Publication Year: **1771** Number: **1**

Printed At: **Edinburgh** Shelf Locator: **EB.1**

Physical Description: **3 v., 160 plates : ill. ; 26 cm. (4to)**

MMSID **992277653804341** Genre: **encyclopedia**

Language: **English** Number of Volumes: **3**

### Volumes for Edition 1,1771

Search Volume

- [2 C-L](#)
- [3 M-Z](#)
- [1 A-B](#)

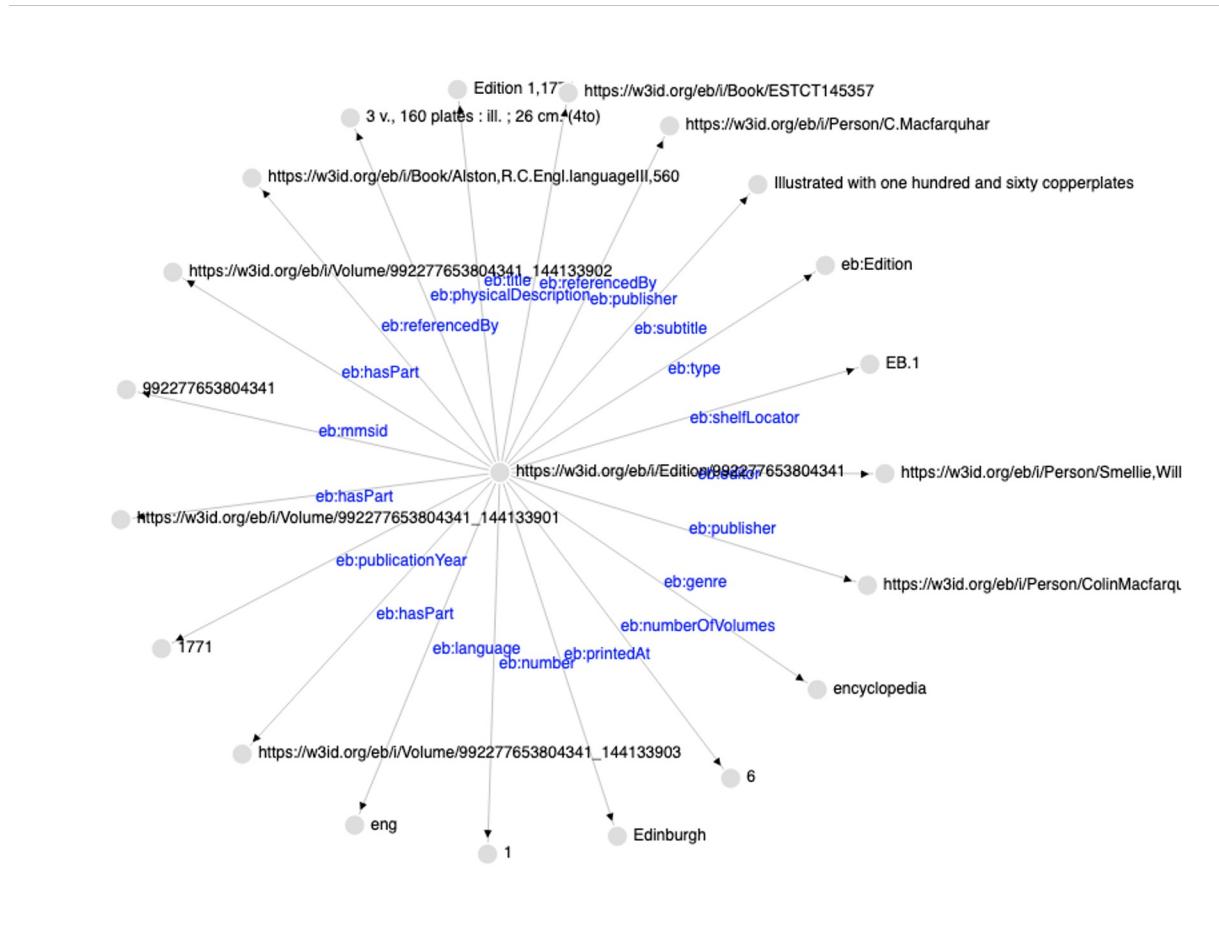
Title: **Encyclopaedia Britannica; or, A dictionary of arts and sciences, compiled upon a new plan**

Permanent URL: <https://digital.nls.uk/144133902>

ID: **144133902** Number: **2**

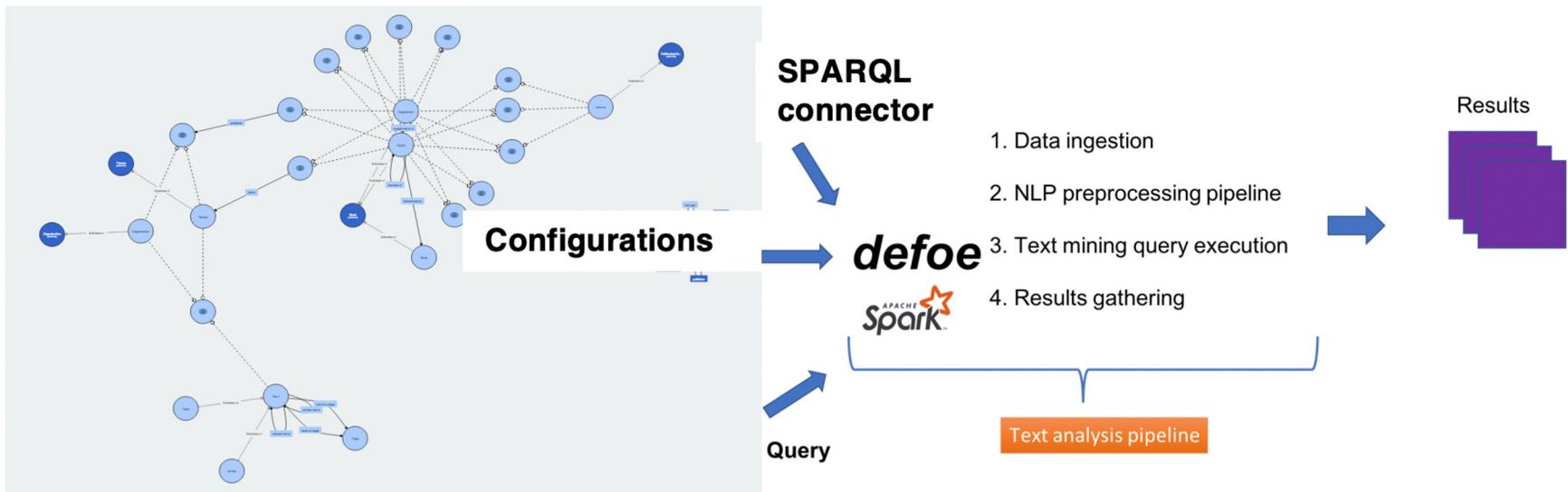
Number of Pages: **1018** Letters: **C-L**

## Visualization of Resources



Visualization of the First Edition (issued in 1771) resource stored in the EB-KG

## Defoe Queries – Text Mining Platform



```
def do_query(archives, config_file=None, logger=None, \
            context=None):
    # [archive, archive, ...]
    documents = archives.flatMap(lambda archive:\n        list(archive))
    # [num_pages, num_pages, ...]
    num_pages = documents.map(lambda document: \
        document.num_pages)
    result = [documents.count(), num_pages.reduce(add)]
    return {"num_volumes": result[0], \
            "num_pages": result[1]}
```

*total\_pages* query: Iterates through archives and counts the total number of documents (e.g. volume, book etc) and total number of pages.



Result:  
num\_volumes: 10  
num\_pages: 9448

Results using the 2<sup>nd</sup> Edition – EB

frances

Term Search   Term Similarity   Topic Modelling   Defoe Queries   Collection Details

DY

## Defoe Queries

Collection  
Encyclopaedia Britannica ▾

Query type  
frequency\_keysearch\_by\_year

Encyclopaedia Britannica

Chapbooks printed in Scotland

Gazetteers of Scotland

Ladies' Edinburgh Debating Society

### HIT Count

It counts the number of times a keyword/keysentence appears in a term (article or topic).

Word  Term

### Upload Lexicon File

The file should contain a line per keyword and/or key sentence that you want to use in your query.

[UPLOAD FILE](#)

### Filtering Options (Optional)

Target words or sentences

Word or Sentence

[ADD](#)

Select term which contains **Any** the target words/sentences

Or contains **All** of the target words/sentences.

Year Range

1768

To

1860

[SUBMIT QUERY](#)

## Defoe Queries – Text Mining Platform

frances

Term Search

Term Similarity

Topic Modelling

Defoe Queries

Collection Details

DY

### Defoe Queries

Collection

Encyclopaedia Britannica ▾

Query type

frequency\_keysearch\_by\_year

It counts the number of terms/words in sentences. It groups results by years.

frequency\_keysearch\_by\_year

publication\_normalized

uris\_keysearch

snippet\_keysearch\_by\_year

fulltext\_keysearch\_by\_year

geoparser\_by\_year

frequency-distribution

lexicon-diversity

person\_entity\_recognition

None ▾

Word  Term

UPLOAD FILE

Preprocess Treatment

It does not apply any type of treatment to the text.

Hit Count

It counts the number of times a keyword/keysentence appears in a term (all words).

Upload Lexicon File

The file should contain a line per keyword and/or key sentence that you want to search for.

Filtering Options (Optional)

Target words or sentences

Word or Sentence

Select term which contains **Any** the target words/sentences

Or contains **All** of the target words/sentences.

Year Range

1768  To  1860

# Defoe Queries – Text Mining Platform

frances

[Term Search](#)
[Term Similarity](#)
[Topic Modelling](#)
[Defoe Queries](#)
[Collection Details](#)

DY

-  Profile
-  Saved
-  Tasks
-  Sign out

## Defoe Query Tasks

Collection	Query Type	Configuration	State	Submit Time	Actions
Chapbooks printed in Scotland	geoparser_by_year	animal.txt	DONE	2023-05-31 20:02:51.880409	<a href="#">VIEW</a>
Encyclopaedia Britannica	frequency_keysearch_by_year	commodities.txt lemmatize term	DONE	2023-06-13 18:41:38.677508	<a href="#">VIEW</a>
Ladies' Edinburgh Debating Society	publication_normalized		DONE	2023-06-13 18:42:10.114119	<a href="#">VIEW</a>
Chapbooks printed in Scotland	geoparser_by_year	scots.txt lemmatize	DONE	2023-06-14 11:55:08.393946	<a href="#">VIEW</a>
Chapbooks printed in Scotland	geoparser_by_year	scots.txt normalize	DONE	2023-06-14 13:10:14.124065	<a href="#">VIEW</a>
Chapbooks printed in Scotland	geoparser_by_year	scots.txt	DONE	2023-06-20 08:35:23.817690	<a href="#">VIEW</a>
Chapbooks printed in Scotland	geoparser_by_year	scots.txt lemmatize	DONE	2023-06-20 08:47:17.244466	<a href="#">VIEW</a>
Encyclopaedia Britannica	publication_normalized		DONE	2023-06-21 20:42:21.304656	<a href="#">VIEW</a>
Chapbooks printed in Scotland	publication_normalized		DONE	2023-06-21 20:42:36.779318	<a href="#">VIEW</a>

## The Journey

Phase 1: Information Extraction

Phase 2: Ontology and Knowledge Graphs

Phase 3: Augmented Knowledge Graphs with Deep Transfer Learning

**Phase 4: Defoe and Knowledge Graphs**

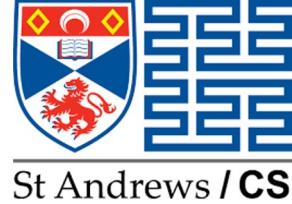
Phase 5: React-Flask Web Platform

**Phase 6: Case Studies**

6.1 Commodity and Enslaved Labor in the EB

6.2 Geoparsing the Chapbooks printed in Scotland

# Phase 6.1: Commodity and Enslaved Labor in the EB



Collaboration with: Ash Charlton & Melissa Terras

This case study examines how race, slavery, and products made by enslaved labor were portrayed in the Encyclopaedia Britannica (EB) during a time of intellectual and industrial growth.

This case study conducts an analysis of references to these topics, with a specific focus on **cotton**.

Collaboration with: Ash Charlton & Melissa Terras

frances

Term Search   Term Similarity   Topic Modelling

**Defoe Queries**

Collection Details

DY

## Defoe Queries

Collection  
Encyclopaedia Britannica ▾

Query type  
frequency\_keysearch\_by\_year

It counts the number of terms/words i

frequency\_keysearch\_by\_year

sentences. It groups results by years.

### Preprocess Treatment

It does not apply any type of treatment to the text.

publication\_normalized

None ▾

### Hit Count

It counts the number of times a keyword/keysentence appears in a term (a

uris\_keysearch

Word  Term

### Upload Lexicon File

The file should contain a line per keyword and/or key sentence that you w

snippet\_keysearch\_by\_year

UPLOAD FILE

### Filtering Options (Optional)

fulltext\_keysearch\_by\_year

Target words or sentences

Word or Sentence

ADD

Select term which contains **Any** the target words/sentences

Or contains **All** of the target words/sentences.

Year Range

1768

To

1860

# Phase 6.1: Commodity and Enslaved Labor in the EB



St Andrews / CS

Collaboration with: Ash Charlton & Melissa Terras

Collection

Encyclopaedia Britannica ▾

Query type

frequency\_keysearch\_by\_year

It counts the number of terms/words in which appear your selected keywords/keysentences. It groups results by years.

## Preprocess Treatment

It normalizes the text first and lemmatizes it later, returning the base (lemma) of each word.

Normalize & Lemmatize ▾

## Hit Count

It counts the number of times a keyword/keysentece appears in a term (article or topic).

Word  Term

## Upload Lexicon File

The file should contain a line per keyword and/or key sentence that you want to use in your query.



UPLOAD FILE

20230723-142924\_commodities.txt

## Filtering Options (Optional)

Target words or sentences

Word or Sentence

ADD

Select term which contains **Any** the target words/sentences

Or contains **All** of the target words/sentences.

Year Range

1768



To

1860

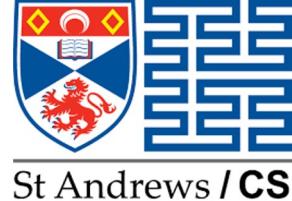


SUBMIT QUERY



Commodities associated with slavery: 'tobacco', 'rum', 'indigo', 'sugar', coffee', 'rice', 'cotton' and 'molasses'

# Phase 6.1: Commodity and Enslaved Labor in the EB



Defoe Queries

Collaboration with: Ash Charlton & Melissa Terras

Collection: Encyclopaedia Britannica

Query Type: frequency\_keysearch\_by\_year

Lexicon Filename: commodities.txt

Preprocess: Normalize & Lemmatize

Hit Count: Term

Year range: 1768 - 1860

Submitted time: 2023-05-16 12:29:34.445029

[CREATE ANOTHER QUERY](#)

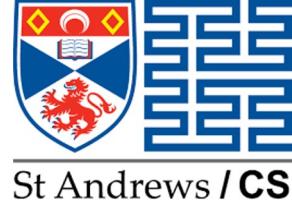
[CHECK ALL QUERY TASKS](#)

Result

[DOWNLOAD](#)

Year	Lexicon	Frequency
1771	sugar	35
	coffee	3
	indigo	13
	rice	12
	cotton	31
	tobacco	15
1773	rum	9
	rice	9
	cotton	33
	indigo	16

# Phase 6.1: Commodity and Enslaved Labor in the EB



Defoe Queries

Collaboration with: Ash Charlton & Melissa Terras

Collection: Encyclopaedia Britannica

Query Type: frequency\_keysearch\_by\_year

Lexicon Filename: commodities.txt

Preprocess: Normalize & Lemmatize

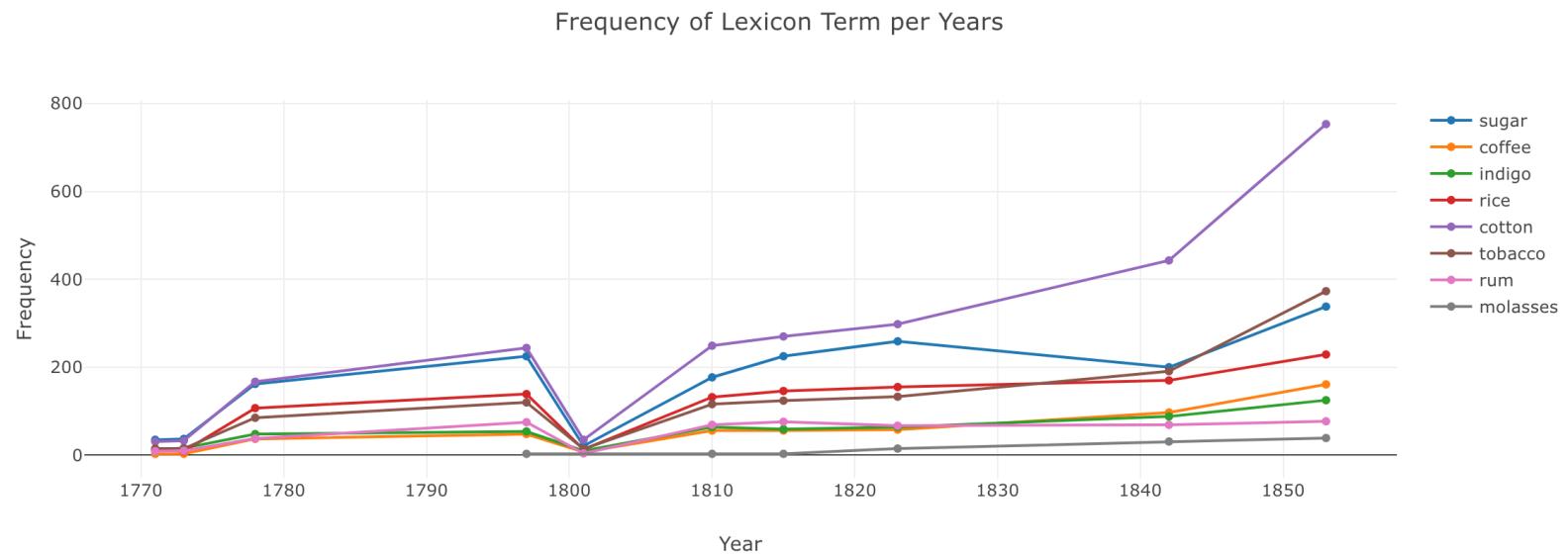
Hit Count: Term

Year range: 1768 - 1860

Submitted time: 2023-05-16 12:29:34.445029

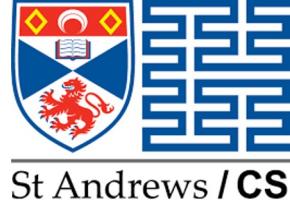
[CREATE ANOTHER QUERY](#)

[CHECK ALL QUERY TASKS](#)

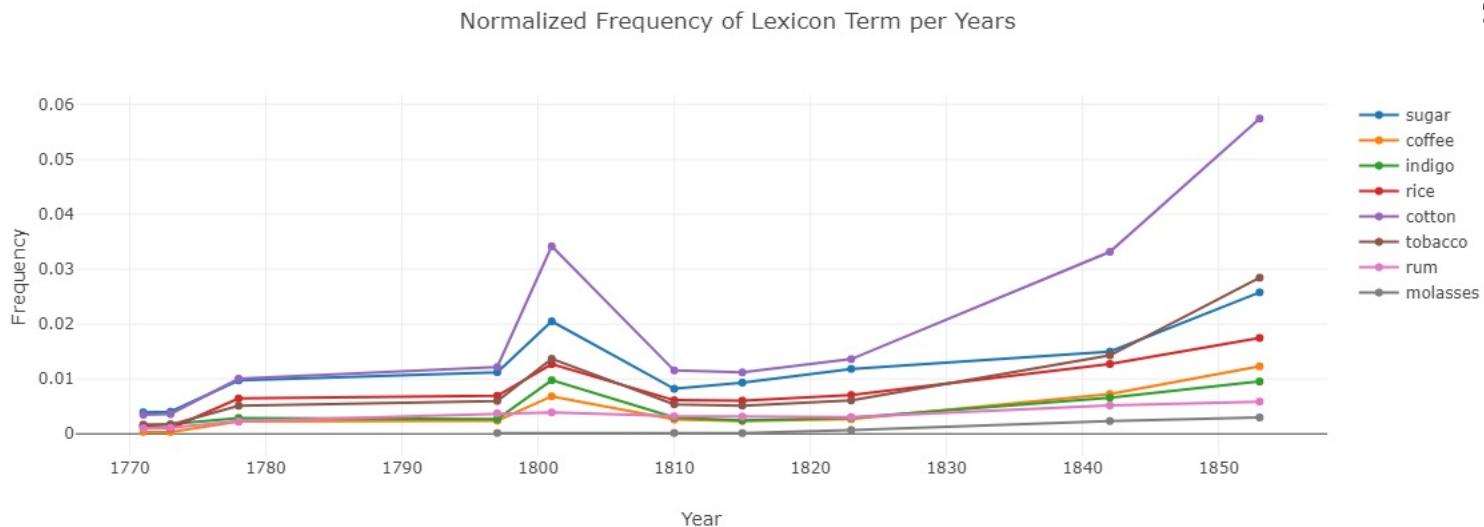


*N-gram to visualise the frequencies of eight commodities*

# Phase 6.1: Commodity and Enslaved Labor in the EB



Collaboration with: Ash Charlton & Melissa Terras



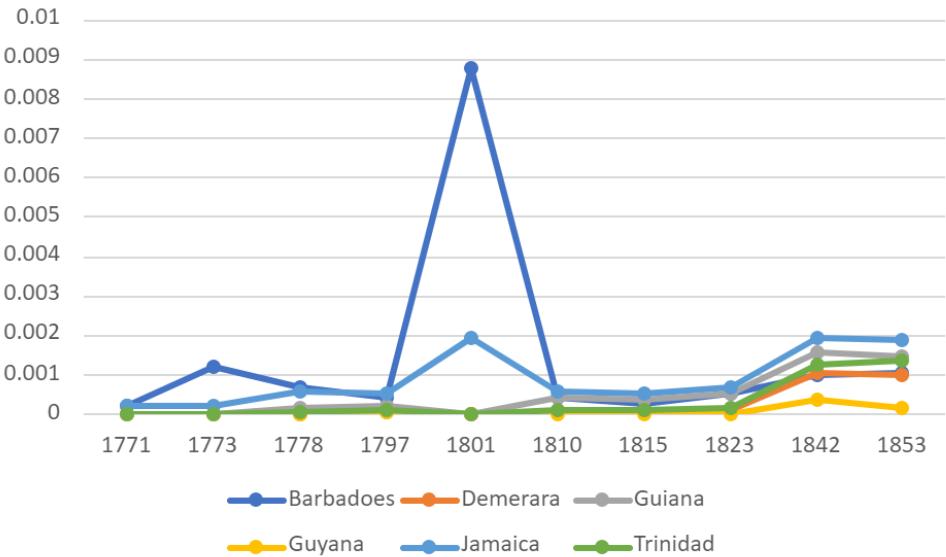
*N-gram to visualise the normalized frequencies of eight commodities*

The analysis uncovered shifts in the frequency of commodity mentions over time, with **'cotton' surpassing 'sugar' as the most referenced commodity after 1801**.

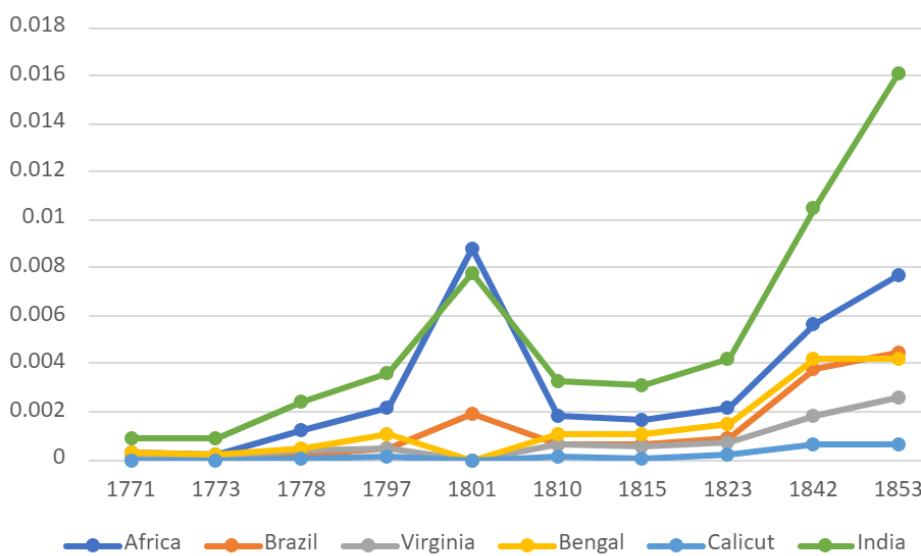
The notable increase in terms mentioning cotton can be attributed to Britain's **advancements in cotton technologies during the late eighteenth and nineteenth centuries**, facilitating mass production and manufacturing

# Phase 6.1: Commodity and Enslaved Labor in the EB

Collaboration with: Ash Charlton & Melissa Terras



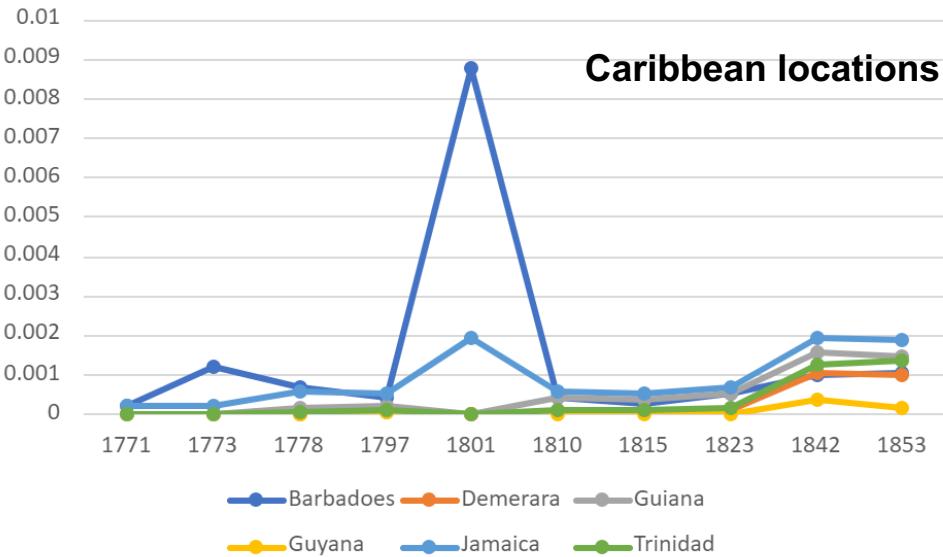
Occurrences of 'cotton' in EB terms mentioning **Caribbean locations**: 'Barbadoes', 'Demerara', 'Guiana', 'Guyana', 'Jamaica' and 'Trinidad'



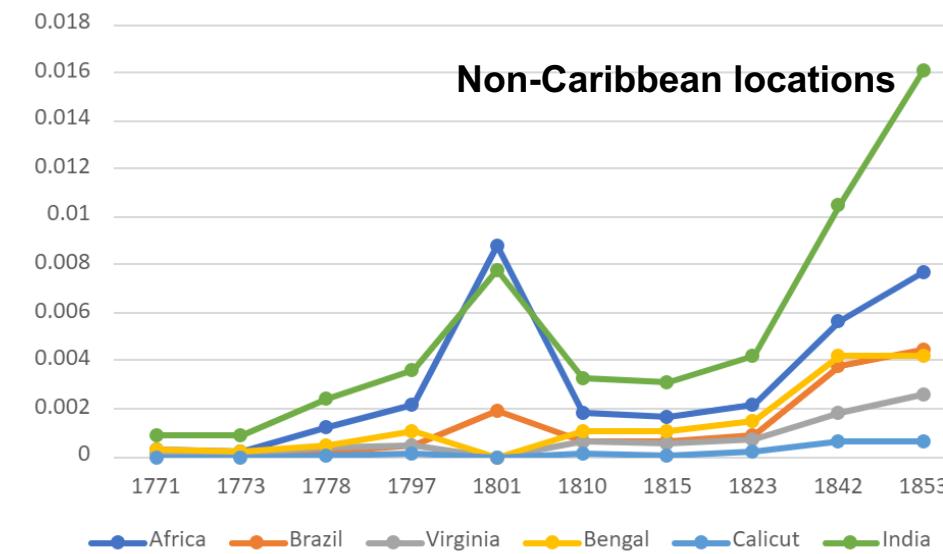
Occurrences of 'cotton' in EB terms mentioning **non-Caribbean** places: 'Africa', 'Brazil', 'Virginia', 'Bengal', 'Calicut', 'India'

# Phase 6.1: Commodity and Enslaved Labor in the EB

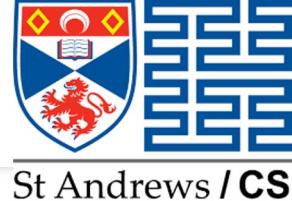
Collaboration with: Ash Charlton & Melissa Terras



The analysis shows a shift in cotton references from the Caribbean to India after 1823, likely due to Britain's East India Company's trading activities.



# Phase 6.2: Geoparsing the Chapbooks printed in Scotland



## Defoe Queries

Collection —

Chapbooks printed in Scotland ▾

Query type —

geoparser\_by\_year ▾

It geo-locates locations in pages and geo-resolves them using the Edinburgh Geoparser. It groups results by years.

**Preprocess Treatment**

It does not apply any type of treatment to the text.

**Upload Lexicon File**

The file should contain a line per keyword and/or key sentence that you want to use in your query.

**Gazetteer**

A world-wide gazetteer of over eight million placenames, made available free of charge.

**None ▾**

- None
- Normalize
- Normalize & Numbers
- Normalize & Lemmatize
- Normalize & Stemming

**Filtering Options (Optional)**

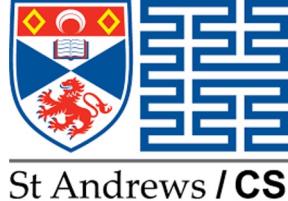
Target words or sentences

Word or Sentence **ADD**

Select term which contains **Any** the target words/sentences

Or contains **All** of the target words/sentences.

# Phase 6.2: Geoparsing the Chapbooks printed in Scotland



Collection — **Chapbooks printed in Scotland**

Query type — **geoparser\_by\_year**

It geo-locates locations in pages and geo-resolves them using the Edinburgh Geoparser. It groups results by years.

**Preprocess Treatment**  
It normalizes the text first and lemmatizes it later, returning the base (lemma) of each word.

**Normalize & Lemmatize**

**Upload Lexicon File**  
The file should contain a line per keyword and/or key sentence that you want to use in your query.

**UPLOAD FILE**  
20230723-143423\_scots.txt

**Gazetteer**  
A world-wide gazetteer of over eight million placenames, made available free of charge.

**Geonames**

**Filtering Options (Optional)**

**Target words or sentences**

**Word or Sentence**

Select term which contains **Any** the target words/sentences  
Or contains **All** of the target words/sentences.

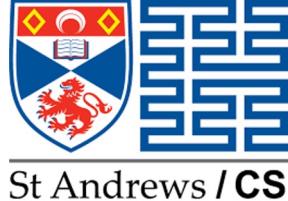
**Year Range**  
**1700**

**Bounding Box**  
The geoparser will prefer places within bounding box, but will still choose locations outside it if other factors give them higher weighting.

**EDIT**

- Geonames
- OS
- Natural Earth
- Geonames through Unlock
- Unlock
- DEEP
- Pleiades+

# Phase 6.2: Geoparsing the Chapbooks printed in Scotland



frances

Term Search Term Similarity Topic Modelling Defoe Queries Collection Details DY

Chapbooks printed in Scotland ▾ geoparser\_by\_year

It geo-locates locations in pages and geo-resolves them using the Edinburgh Geoparser. It groups results by years.

Preprocess Text

Upload Lexicon

Gazetteer

Filtering Options

Target words

Select term variants

Year Range

Bounding Box

Chapbooks printed in Scotland

geoparser\_by\_year

It geo-locates locations in pages and geo-resolves them using the Edinburgh Geoparser. It groups results by years.

Edit bounding box

You can either manually enter the bounding box or search and select a place. The format of bounding box is: W N E S, where W(est) N(orth) E(ast) S(outh) are decimal degrees

West: -14.015517 North: 54.4339831 East: -0.3209221 South: 61.061

Search a place for its bounding box: Sco

Scotland, United Kingdom

Aktau Airport, KR-11, Tüpqarağan District, Mangystau Region, Kazakhstan

Super-Collecteur Ouest, arrondissement de Hay Hassani, مقاطعة الحي الحسني, préfecture d'arrondissement de Hay Hassani, Pachalik de Casablanca, Prefecture of Casablanca, Casablanca-Settat, Morocco

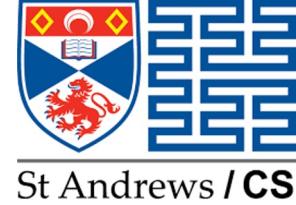
Super-Collecteur Ouest, arrondissement d'Anfa, مقاطعة أنفا, préfecture d'arrondissements de Casablanca-Anfa, عمالة أنفا, Pachalik de Casablanca, Prefecture of Casablanca, Casablanca-Settat, 20202, Morocco

São Cristóvão, Rua Bartolomeu de Gusmão, São Cristóvão, Zona Central do Rio de Janeiro, Rio de Janeiro, Região

CLOSE CONFIRM EDIT

The geoparser will prefer places within bounding box, but will still choose

## Phase 6.2: Geoparsing the Chapbooks printed in Scotland



# Defoe Queries

Collection: **Chapbooks printed in Scotland**

Query Type: geoparser\_by\_year

Lexicon Filename: scots.txt

Preprocess: Normalize & Lemmatize

## Gazetteer: Geonames

**Year range:** 1700 - 1899

Bounding Box: -14.015517 54.4339831 -0.3209221  
61.061

Submitted time: 2023-06-14 11:55:08.393946

[CREATE ANOTHER QUERY](#)

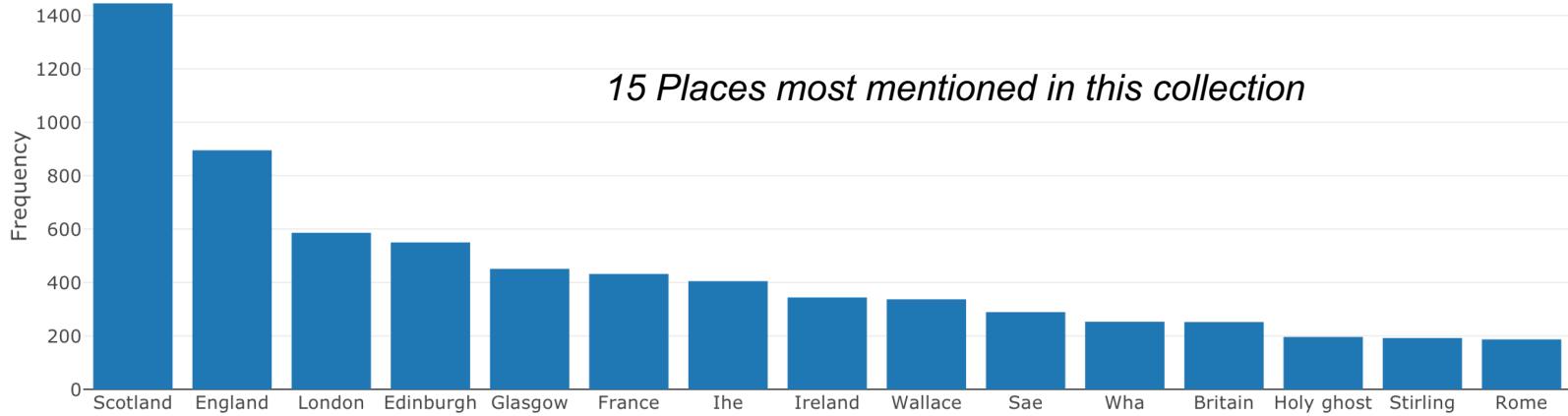
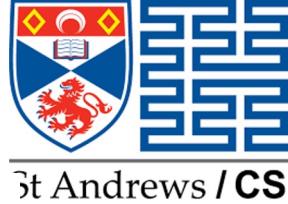
## CHECK ALL QUERY TASKS

Result

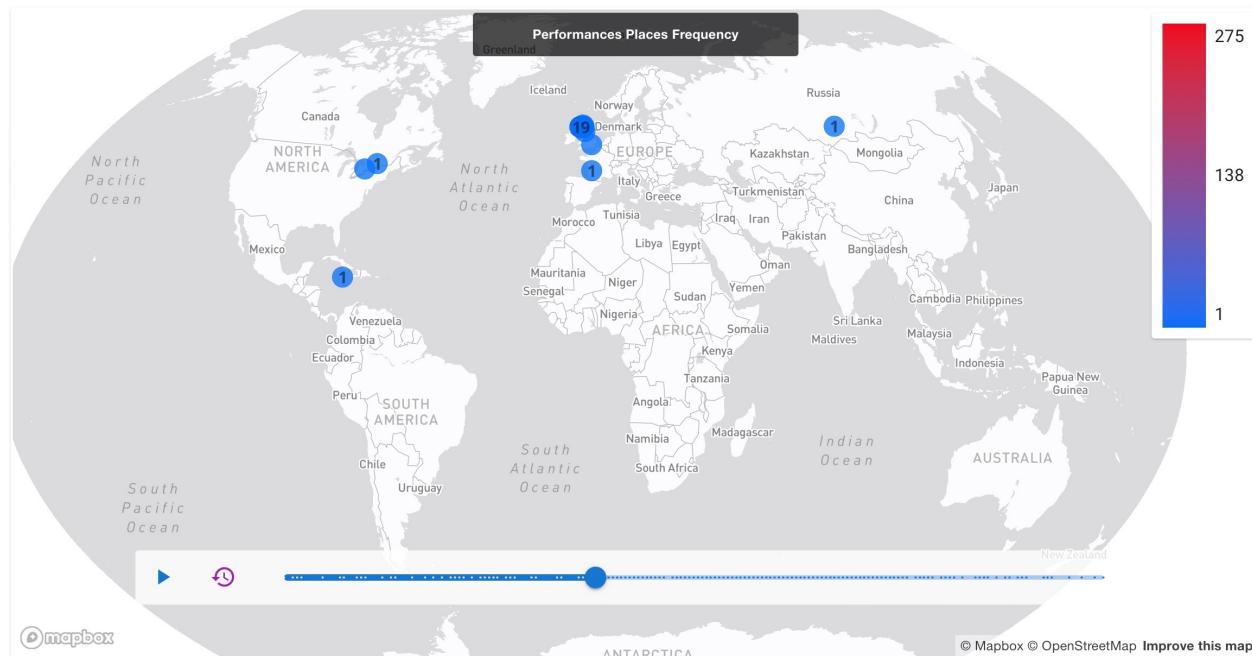
 DOWNLOAD

Year	Series	Volume	Volume ID	Volume Title	Page	Words	Part	Geo
1701	204	1	104184313	unnatural sic son	5	267	None	Spi
	204	1	104184313	unnnatural sic son	3	274	None	Tryal
	204	1	104184313	unnnatural sic son	7	276	None	Sirj
	204	1	104184313	unnnatural sic son	1	170	None	Tryal
	204	1	104184313	unnnatural sic son	6	269	None	Man the
	204	1	104184313	unnnatural sic son	2	280	None	Debaucherie CO Wi
	206	1	104184316	King of France His catechism	3	167	None	Europe
1703	206	1	104184316	King of France His catechism	6	354	None	Italy Vigo Wha

# Phase 6.2: Geoparsing the Chapbooks printed in Scotland

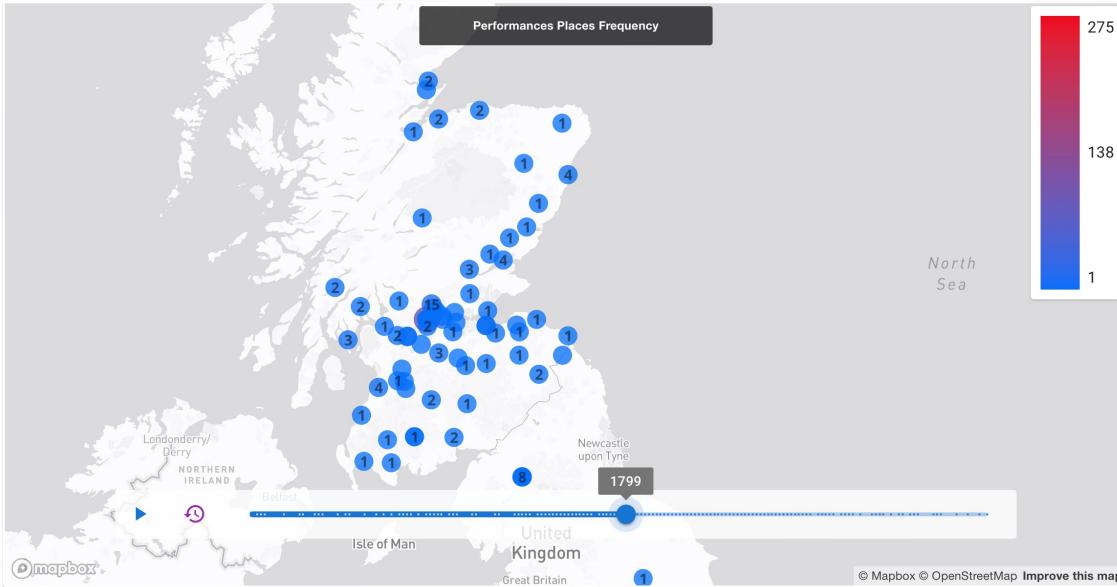


Map of geo-resolved places frequencies over time



By geoparsing the chapbooks using *frances*, historians gain the ability to unravel the spatial dimensions of Scottish literature from the 18th and 19th centuries.

# Phase 6.2: Geoparsing the Chapbooks printed in Scotland



*Zoom in ‘Scotland’.  
Frequency of geolocated places in 1799*



*Geolocation of “Edinburgh”*

# *frances*: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale

## Conclusions

- *frances* enables for analysing digital historical textual collections - scalability and performance
- Cloud-based architecture ensures maintainability, efficient deployment, and scalability
- NLS and EB Ontologies standardizes representation of NLS Data Foundry collections
- Enabling knowledge graphs creation
- Deep Learning Analysis - New suite of ML functionalities that can be used to analyse textual collect.
- Parallel Text Mining queries

## Future Work

- Expanding knowledge graph creation, collaborating with other institutions, and extensive evaluation
- Integrating external linked data sources, and gathering user feedback

## Publications

[\*\*frances: a deep learning NLP and text mining web tool to unlock historical digital collections: a case study on the Encyclopaedia Britannica\*\*](#) Filgueira, R., 14 Dec 2022, 2022 IEEE 18th International Conference on e-Science (e-Science)

[\*\*frances: cloud-based historical text mining with deep learning and parallel processing\*\*](#) Yu, L., Charlton, A., Askins, W., Terras, M. & [Filgueira, R.](#), 12 Jul 2023, (Accepted/In press) 2023 IEEE 19th International Conference on e-Science (e-Science)

# Encyclopædia Britannica;

James OR, A Fullerton

## DICTIONARY

O F

A R T S and S C I E N C E S,

COMPILED UPON A NEW PLAN.

I N W H I C H

The diferent SCIENCES and ARTS are digested into  
distinct Treatises or Systems;

A N D

The various TECHNICAL TERMS, &c. are explained as they occur  
in the order of the Alphabet.

ILLUSTRATED WITH ONE HUNDRED AND SIXTY COPPER

By a SOCIETY of GENTLEMEN in SCOTLA

I N T H R E E V O L U M E S

V O L. I.

E D I N B U R

Printed for A. BELL and C.  
And sold by COLIN MACFARQUHAR, at

M.DCC

# Questions ?

*Frances: A Deep Learning NLP and Text Mining Digital Platform for Analysis of Historical Texts at Scale*

Dr. Rosa Filgueira,

Lecturer at the School of Computer Science,  
University of St Andrews,

Email: rf208@st-Andrews.ac.uk