



Analysing Digital Historical Textual Data at Scale with Defoe Toolbox

Dr. Rosa Filgueira, EPCC,
University of Edinburgh
Email: rosa.filgueira@ed.ac.uk

Context

Working with

- Historians, Humanities and computational linguistics researchers
- Large digital collections been available for research

Funding - March 2018

- ATI-SE Data Engineering Programme
- Living with Machines (LwM)
- Text Data Mining (TDM)
- CDCS Text Data Mining Lab

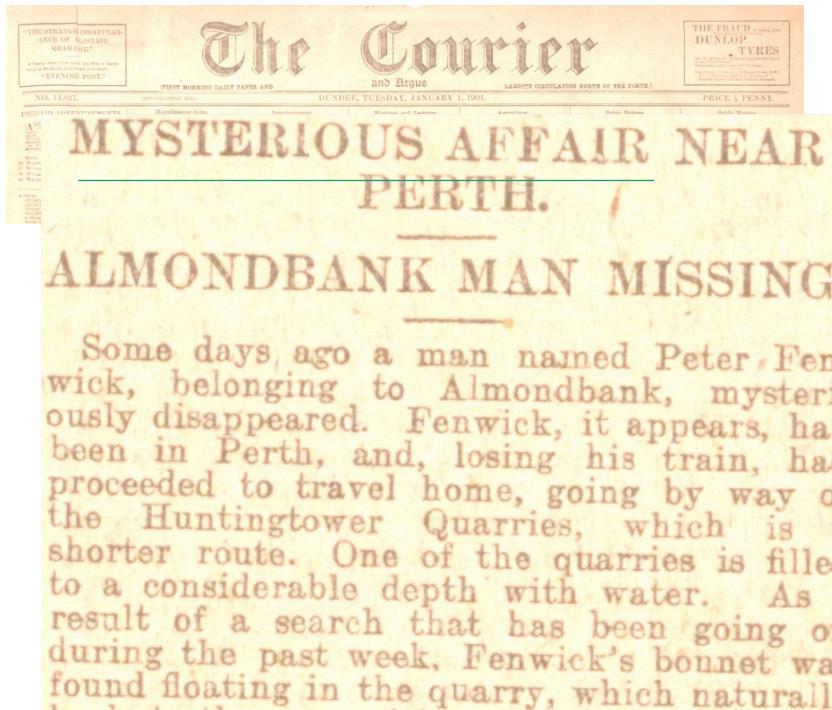
Motivation – eScience for Historians & Humanities communities

- Hunger for large scale text mining facilities
- Limited capacity and/or skills to use:
 - HPC/Cloud environments
 - analytic frameworks to create applications

Context

Challenges

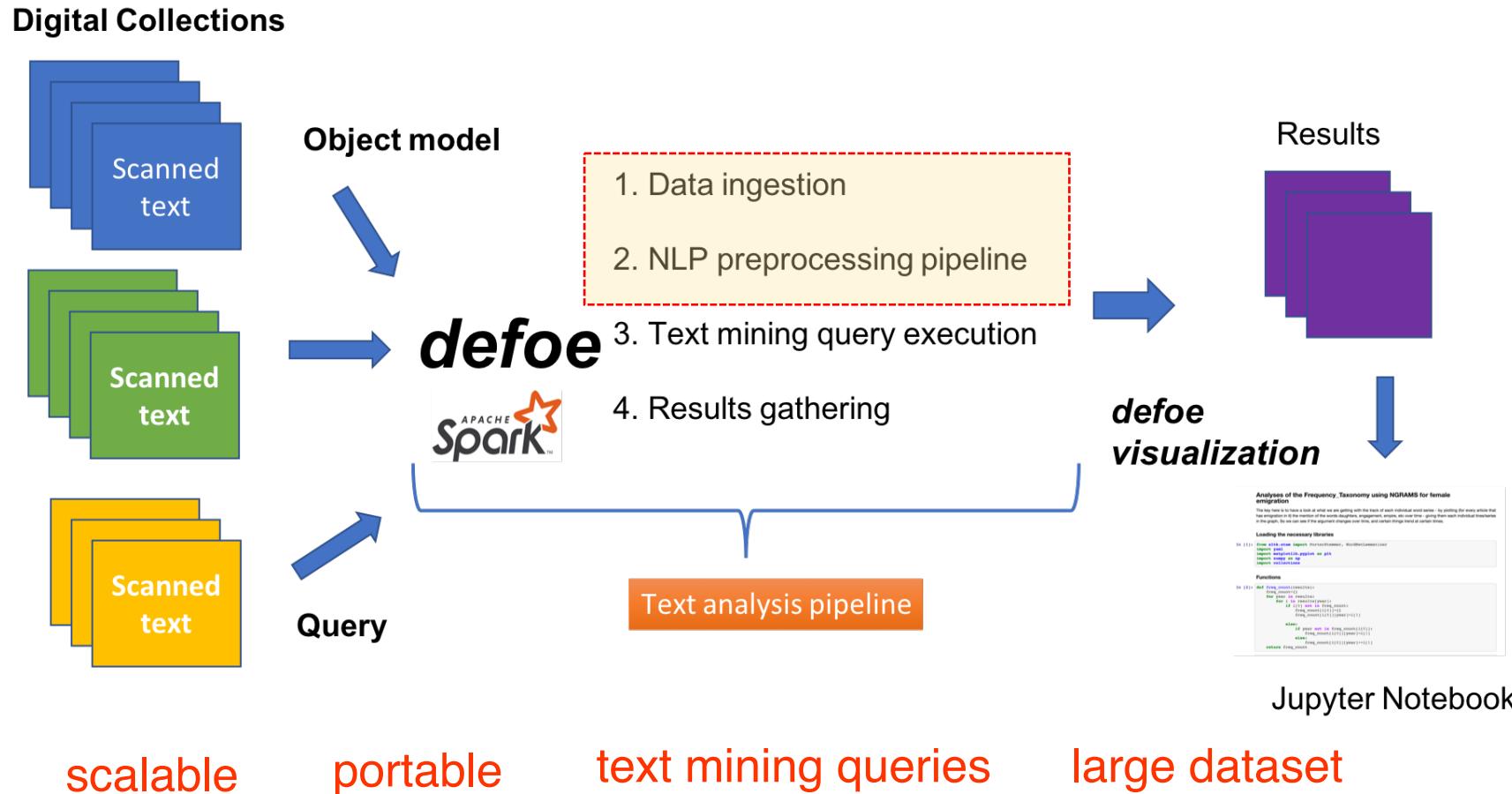
- Several large digital collections (semi-structured data)
 - Different levels of quality of data – OCR
 - Data with different physical representations and schemas



```
...
<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>
  ...

```

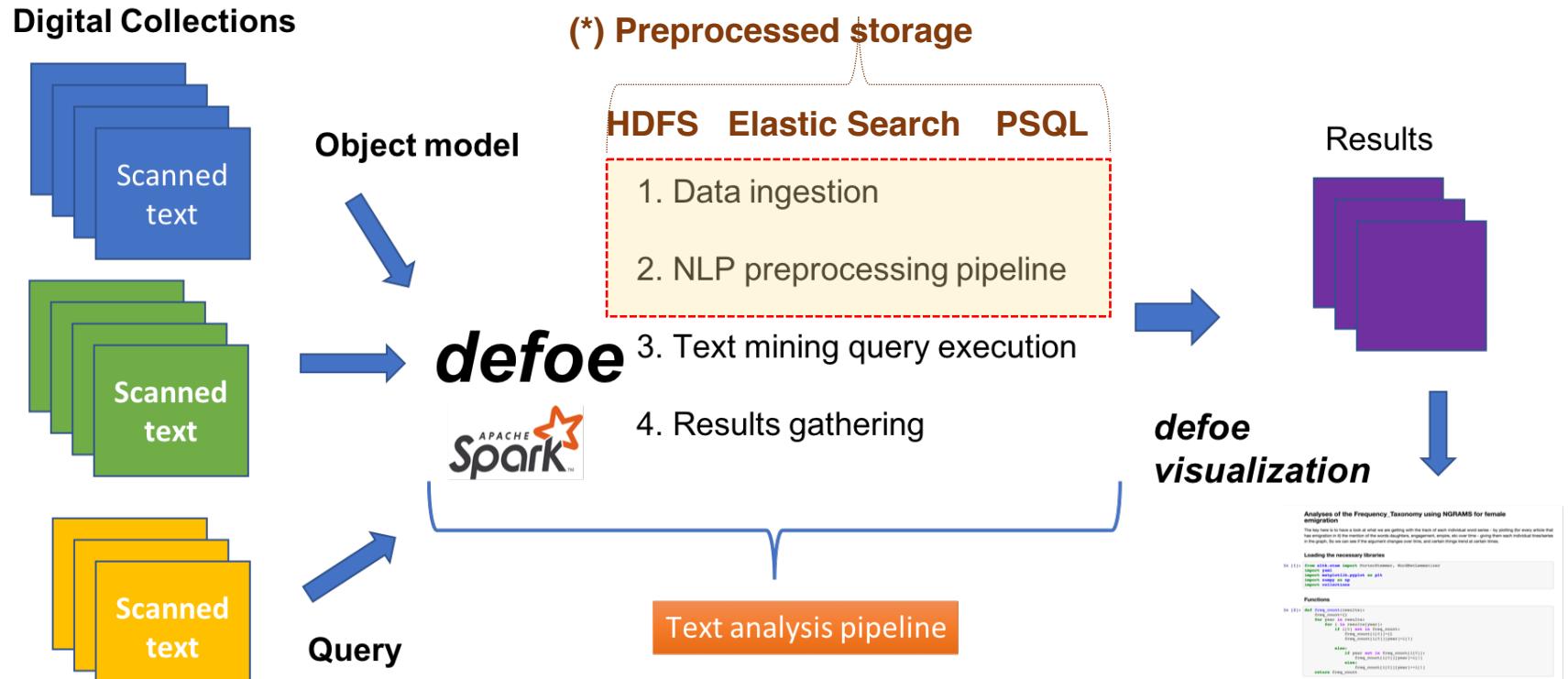
defoe: new eScience toolbox for historical research



<https://github.com/alan-turing-institute/defoe>

https://github.com/alan-turing-institute/defoe_visualization

defoe: new eScience toolbox for historical research



scalable

portable

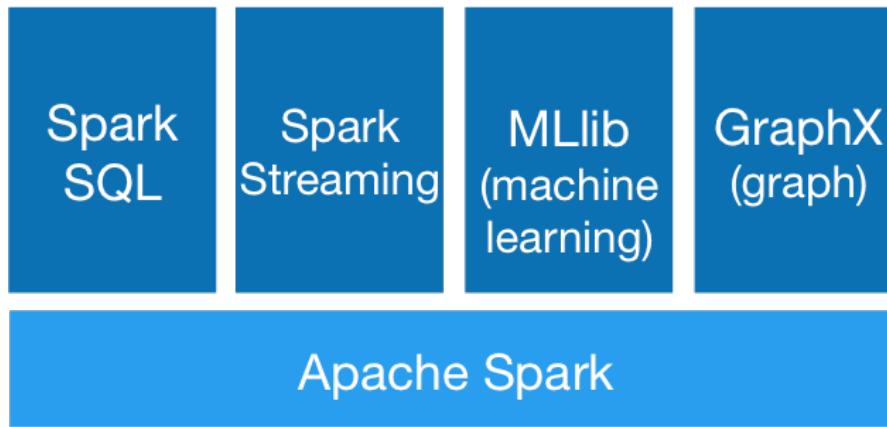
text mining queries

large dataset

<https://github.com/alan-turing-institute/defoe>

https://github.com/alan-turing-institute/defoe_visualization

Apache Spark



- Analytics engine for large-scale data processing
- High performance for batch and streaming data
- APIs
Java, Scala, **Python** and R

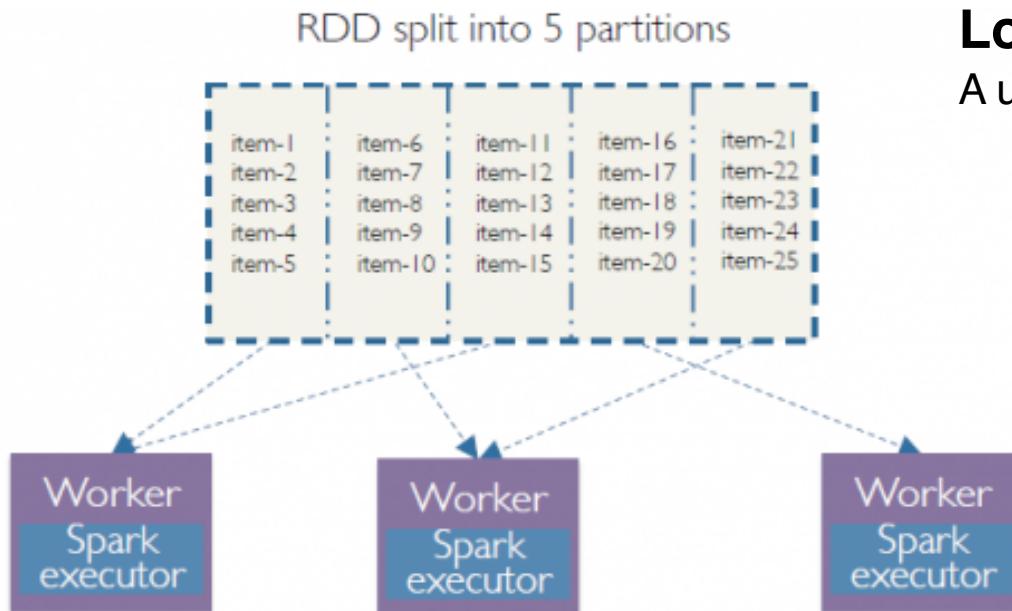
<https://spark.apache.org/>

<https://github.com/EPCCed/prace-spark-for-data-scientists>

Apache Spark

Resilient Distributed Datasets (RDD)

- Represent data or transformations on data
- It is distributed collection of items – partitions
- Read-only → they are immutable
- Enables operations to be performed in parallel



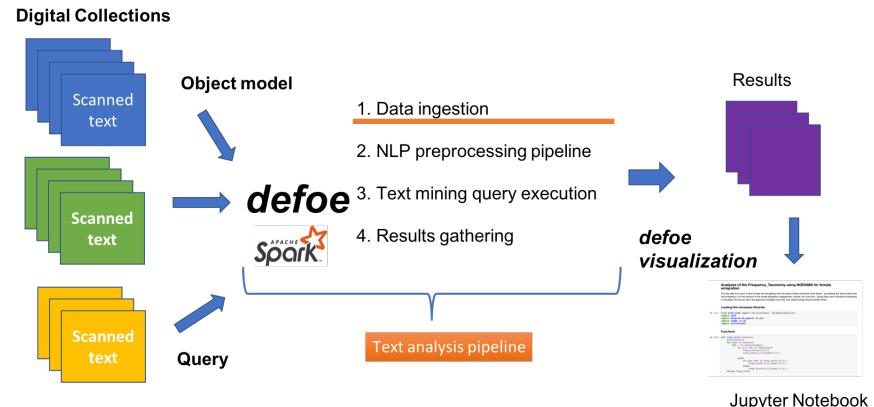
Logical – view

A user sees an object with 25 items

Physically – view

25 items distributed across different CPUS/Cores

Data Ingestion



Support for three physical representations:

- 1) one XML document per issue
- 2) one XML document with search results including several articles
- 3) one XML METS document and ALTO XML per page

Object models -- **loading/parsing** data into RDD:

PAPER

NZPP

ALTO

FMP

NLS

PAPER object model (British Library Newspapers)

RDD



0000164- The Courier and Argus
0000187- The Bath Chronicle
0000195- Archer Bath Chronicle
0000321- The Nottingham Evening Post
0000452- Edinburgh Evening News

Class Issue → Representation of an issue (XML document)

Each XML holds articles that belong to the same issue

issue

filename
issue_tree
issue_id
date
page_count
day_of_the_week

attributes

article list



Class Article → Representation of an article

article

attributes

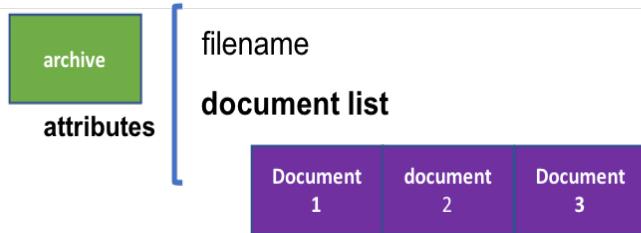
article_tree
filename
quality
title
preamble
content
article_id
page_ids
words = content + title + preamble

ALTO/NLS object models (British Library books and NLS Digital Collection)

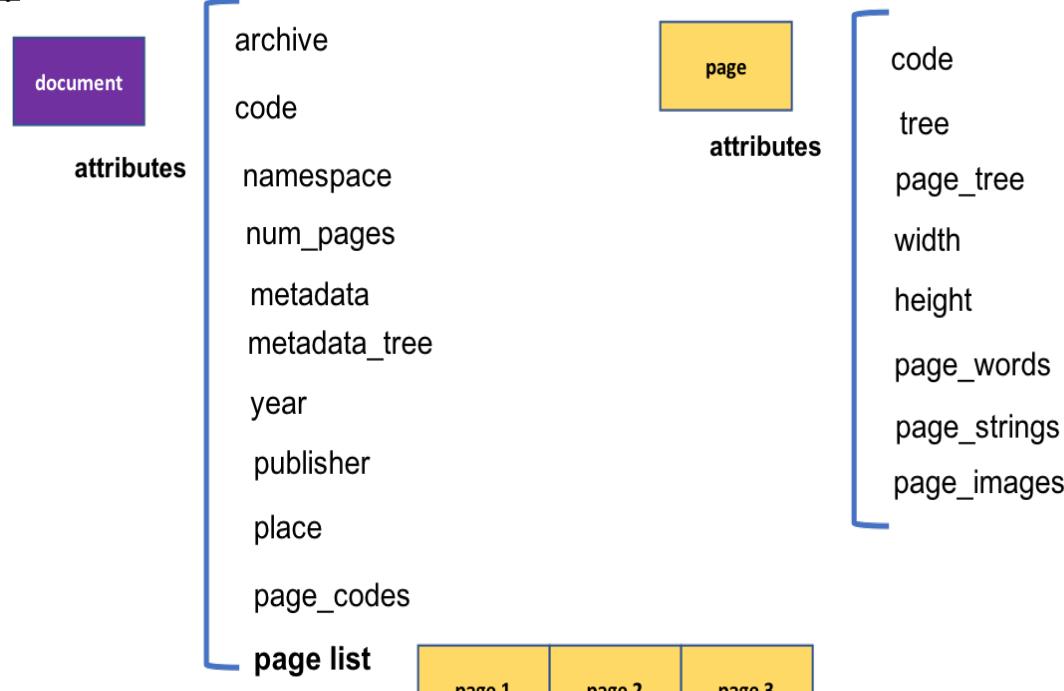


1510_1699/000001143_0_1-20pgs_560409_dat.zip
1510_1699/000000874_0_1-22pgs_570785_dat.zip
1510_1699/000051983_0_1-92pgs_568584_dat.zip
1510_1699/000987728_0_1-92pgs_567840_dat.zip

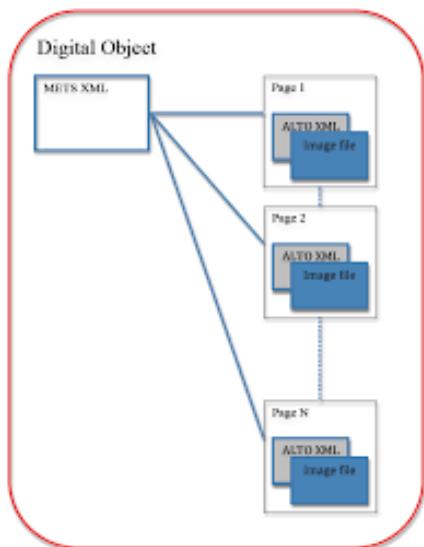
Class Archive → Representation of a zipped archive



Class Document → Representation of a metadata document (XML in METS/MODS) and an ALTO directory



Class Page → Representation of a page (XML in ALTO).



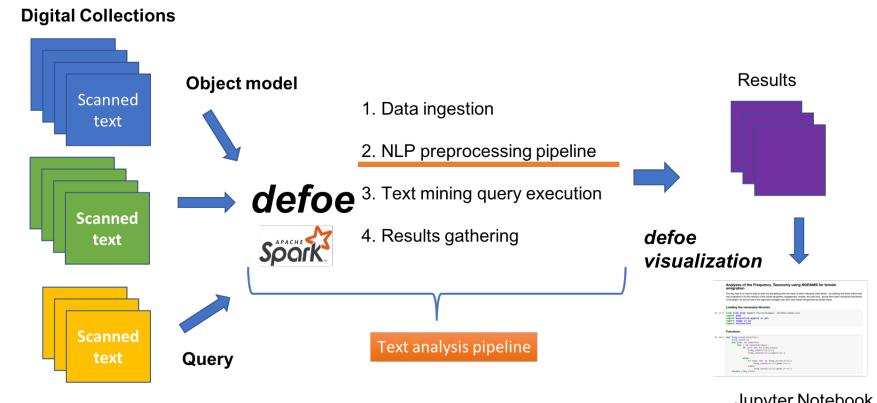
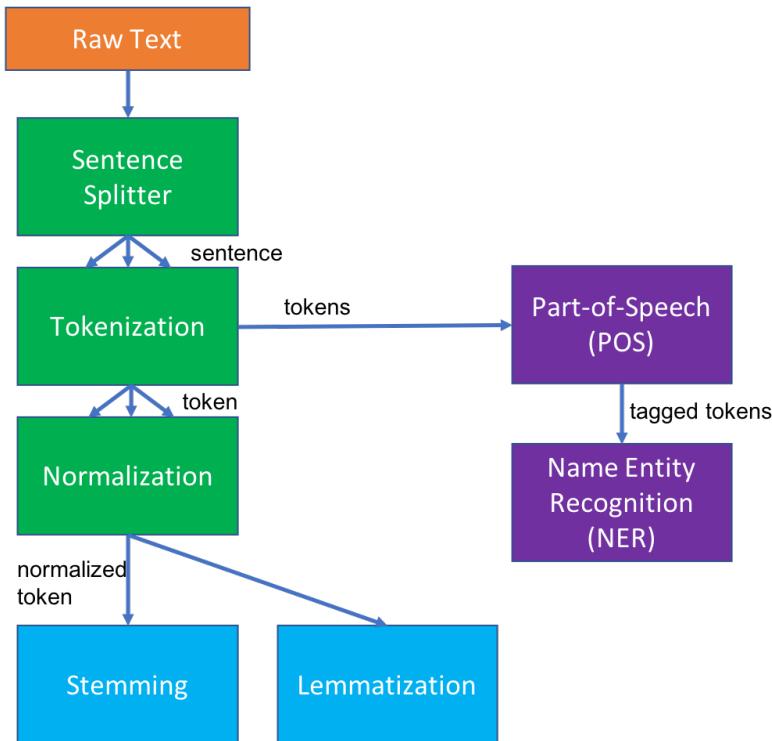
Digital Collections

Dataset	Period	Structure	XML Schema	Space	Model
British Library Books (BLB)	1510-1899	ZIP per book - XML metadata - XML per page	METS and ALTO schemas	~220GB	<u>ALTO</u>
British Library Newspapers (BLN)	1714-1950	XML per issue	GALEN Schema	~1TB	<u>PAPERS</u>
Times Digital Archive (TDA)	1785-2009	XML per issue	GALEN Schema	~324GB	<u>PAPERS</u>
Papers Past New Zealand and Pacific newspapers (NZPP)	1839-1863	XML per 22 articles	XML from a search via an API	~4GB	<u>NZPP</u>
Gazetteers of Scotland, 1803-1901 (*)	1803-1901	- XML metadata - XML per page	METS and ALTO schema	4.7 GB (includ. Fig)	<u>NLS</u>
Encyclopaedia Britannica (**)	1768-1860	- XML metadata - XML per page	METS and ALTO schemas	46GB (includ. figs)	<u>NLS</u>

(*) <https://data.nls.uk/data/digitised-collections/gazetteers-of-scotland/>

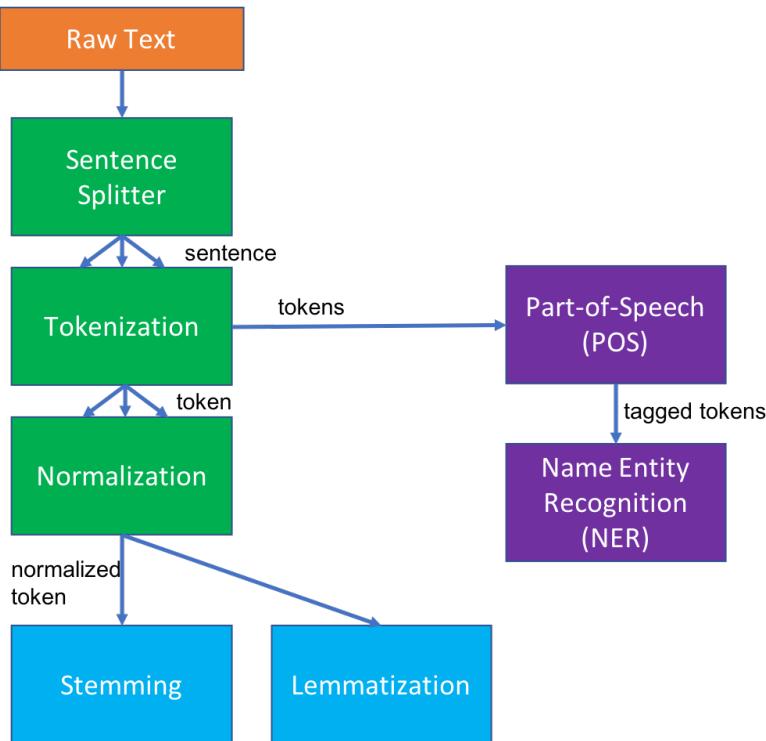
(**) <https://data.nls.uk/data/digitised-collections/encyclopaedia-britannica/>

NLP Preprocessing



Jupyter Notebook

NLP Preprocessing



1. Divide the text in sentences
2. Each word of a sentence is tokenised
3. Each token is normalised
- 4.a) Stemming: reduces words to their roots
- 4.b) Lemmatisation: reduces words to a common base
5. Tag each token
(e.g. nouns, verbs, adjectives, etc.)
6. Classify tagged tokens
(e.g. names of persons, locations, etc)

Spacy and NLTK

Text Mining Queries

NLS Model

```
📄 colocates_by_year.py  
📄 geoparser_pages.py  
📄 georesolution_pages.py  
📄 inventory_per_year.py  
📄 keysearch_by_year.py  
📄 keysentence_by_year_paper.py  
📄 keyword_by_word.py  
📄 keyword_by_year.py  
📄 keyword_concordance_by_word.py  
📄 keyword_concordance_by_year.py  
📄 keyword_metadata_by_word.py  
📄 normalize.py  
📄 ocr_quality_by_year.py  
📄 ocr_quality_multi_level_by_year.py  
📄 preprocessing_sentences.py  
📄 total_documents.py  
📄 total_pages.py  
📄 total_words.py  
📄 write_pages_df_es.py  
📄 write_pages_df_hdfs.py  
📄 write_pages_df_psql.py
```

ALTO Model

```
📄 colocates_by_year.py  
📄 keyword_by_word.py  
📄 keyword_by_year.py  
📄 keyword_concordance_by_word.py  
📄 keyword_concordance_by_year.py  
📄 normalize.py  
📄 ocr_quality_by_year.py  
📄 total_documents.py  
📄 total_pages.py  
📄 total_words.py
```

PAPERS Model

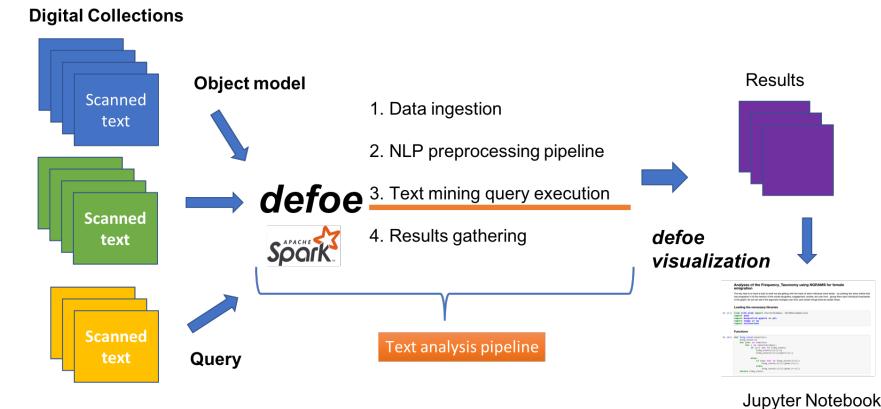
```
📄 colocates_by_year.py  
📄 keysentence_by_year.py  
📄 keyword_by_year.py  
📄 keyword_concordance_by_date.py  
📄 keywords_by_year.py  
📄 lda_topics.py  
📄 normalize.py  
📄 ocr_quality_by_year.py  
📄 target_and_keywords_by_year.py  
📄 target_and_keywords_count_by_ye...  
📄 target_concordance_collocation_b...  
📄 total_articles.py  
📄 total_issues.py  
📄 total_words.py  
📄 unique_words.py
```

NZPP Model

```
📄 keyword_by_year.py  
📄 keyword_concordance_by_date.py  
📄 normalize.py  
📄 total_articles.py  
📄 total_words.py
```

FMP Model

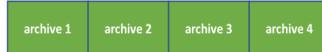
```
📄 keyword_metadata_by_word.py  
📄 keyword_segmentation.py  
📄 normalize.py  
📄 target_segmentation.py  
📄 total_articles.py  
📄 total_documents.py
```



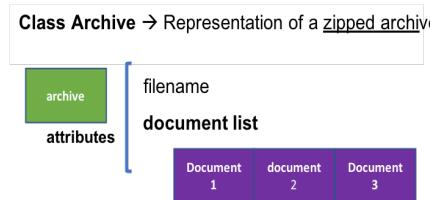
Jupyter Notebook

Text Mining Queries

RDD



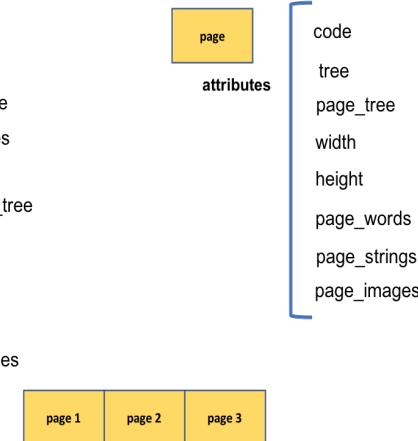
ALTO object model



Class Document → Representation of a metadata document (XML in METS/MODS) and an ALTO directory



Class Page → Representation of a page (XML in ALTO)



total_words query:

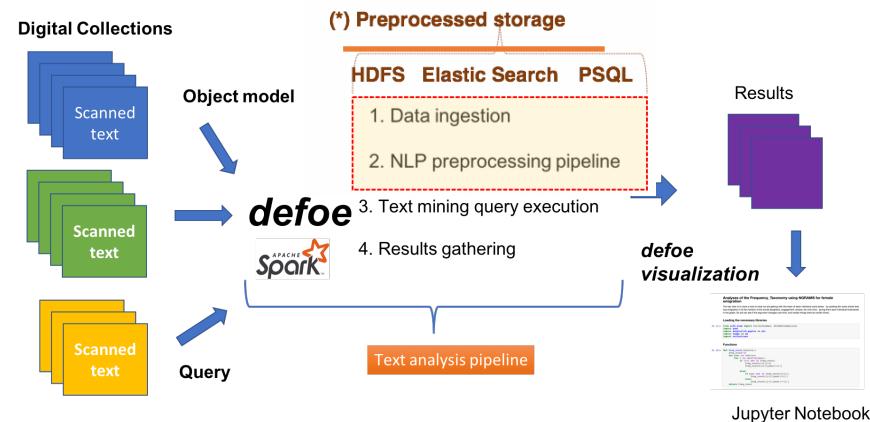
- Iterates through archives
- Count **total number of documents (books)** and **total number of words** among them.

Sample results

Query over British Library Books

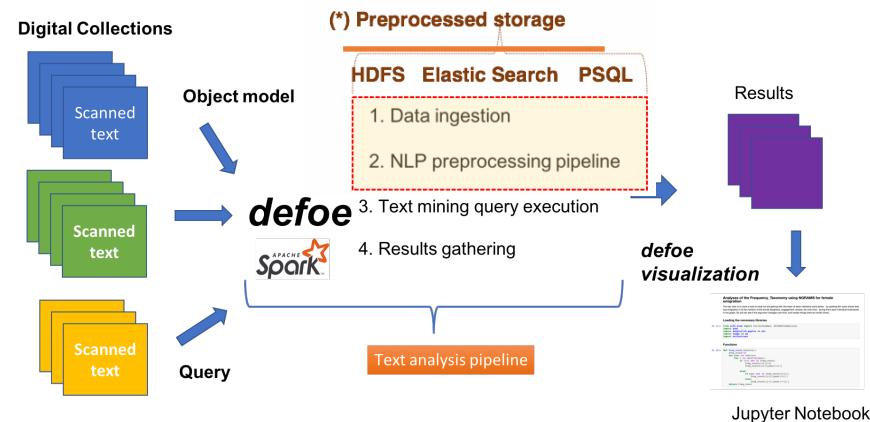
```
{num_documents: 63701, num_words: 6866559285}
```

Preprocessed Data



1. Transform collections to a common format after being preprocessed
2. Run the same queries across collections
3. Capture of metadata + Raw Text + Preprocessed text

Preprocessed Data



Title: Title of the collection - e.g. *Encyclopaedia Britannica*

Edition: Edition of the book/issue - e.g. *Seventh edition, Volume 13, LAB-Magnetism*

Year: Year of publication/edition - e.g. 1842

Place: Place - e.g. *Edinburgh*

archive_filename: Directory Path of the book/newspaper - e.g. */lustre/home/ .../.../193108323*

source_text_filename: Directory Path of the page/issue - e.g. *alto/193201316.34.xml*

type_archive: Type of archive *book / newspaper*

text_unit: Unit that represent each ALTO XML - e.g *Page/Article*

num_text_unit: Number of pages of a book / Number of articles of an issue - e.g. 810

text_unit_id: Id of the page - e.g. *Page3*

model: defoe model used for ingesting this dataset – e.g *nls*

source_text_raw: Text of a page/article

source_text_clean: Cleaned text of page/article – long-S and hyphenated word fixes

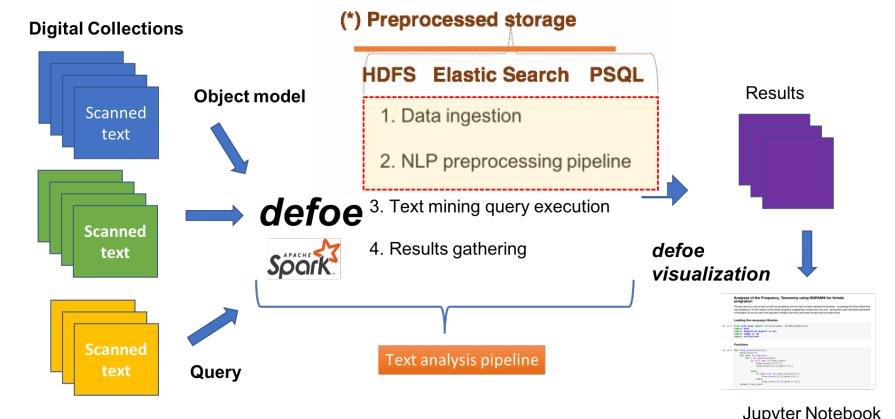
source_text_norm: Normalised cleaned text

source_text_lemmatized: Lemmatized normalized cleaned text

source_text_stemm: stemmed normalized cleaned

num_words: Number of words of a page/article

Preprocessed Data



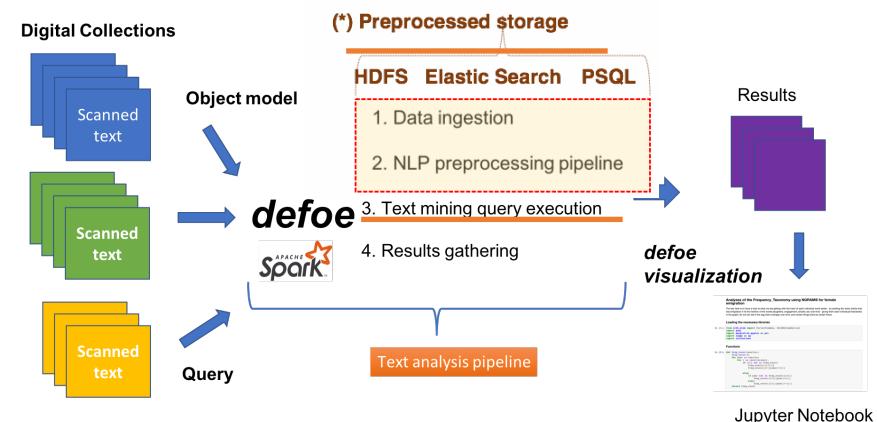
Jupyter Notebook

Example - Scottish Gazetteers

title	edition	year	place	archive_filer	source_text	text_unit	text_unit_id	num_text	type	archive	model	source_text_raw	source_text_clean
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973434/	page	Page1	606	book	nls			iA<* R&o:;,-'.-' ,Äc,Äc- ;<,Äc\'''	i A<* R&o:;,-'.-' ,Äc,Äc- ;<,Äc\'''
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973434/	page	Page2	606	book	nls			^L o,	^L o,
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973434/	page	Page3	606	book	nls			Digitized by the Internet Archive in 2012 http://archive.org/details/	Digitized by the Internet Archive in 2012 http://archive.org/details/
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973434/	page	Page4	606	book	nls			Jac:7T. ./ . 2 1 ttttfUvCGZI cu&cvtc J	Jac:7T. ./ . 2 1 ttttfUv CG ZI cu&cvtc J
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973434/	page	Page5	606	book	nls				
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page6	606	book	nls				
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page7	606	book	nls				
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page8	606	book	nls				
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page9	606	book	nls			THEI GAZETTEER, OF SCOTLAND; CONTAINING A PARTICULAR ANTI TH EI GAZETTEER, OF SCOTLAND; CONTAINING	THEI GAZETTEER, OF SCOTLAND; CONTAINING A PARTICULAR ANTI TH EI GAZETTEER, OF SCOTLAND; CONTAINING
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page10	606	book	nls				
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page11	606	book	nls			INTRODUCTION. GENERAL DESCRIPTION. SGOTLAND, or that part	INTRODUCTION. GENERAL DESCRIPTION. SGOTLAND, or that part
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page12	606	book	nls			iv INTRODUCTION. great opening is termed the Moray Frith, and e	iv INTRODUCTION. great opening is termed the
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page13	606	book	nls			INTRODUCTION. v acres of cultivated, and 14,218,224 acres of unc \' fays a late author	INTRODUCTION. v acres of cultivated, and 14,218,224 acres of unc \' fays a late author
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973435/	page	Page14	606	book	nls			WV INTRODUCTION. Lammermuir hills, in Berwickshire ; the great r	WV INTRODUCTION. Lammermuir hills, in Berwickshire ; the great r
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page15	606	book	nls			INTRODUCTION. vii rising in the mountainous distriet of Lochaber,	INTRODUCTION. vii rising in the mountainous
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page16	606	book	nls			tin INTRODUCTION. the Atlantic ocean. Clouds of the fea are fraug	tin INTRODUCTION. the Atlantic ocean. Clouds of the fea are fraug
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page17	606	book	nls			INTRODUCTION. Is Highland 'Society of Scotland, who offer premii \' fays this almost enthuiaistic writer	INTRODUCTION. Is Highland 'Society of Scotland, who offer premii \' fays this almost enthuiaistic writer
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page18	606	book	nls			y- introduction: End of England, until it were heard, and underftoo	y- introduction: End of England, until it were heard, and underftoo
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page19	606	book	nls			INTRODUCTION. 4 fiifnes, containing coins of Scottifh gold, were f	INTRODUCTION. 4 fiifnes, containing coins of S
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page20	606	book	nls			xii INTRODUCTION. forme of them of great beauty and value. Chalc	xii INTRODUCTION. forme of them of great beauty and value. Chalc
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page21	606	book	nls			xiii land, we may mention the colley, or true Shep	xiii land, we may mention the colley, or true Shep
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973436/	page	Page22	606	book	nls			teem with abundance of trout	teem with abundance of trout
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page23	606	book	nls			%ly INTRODUCTION* RELIGION. It is generally believed, upon the : %ly INTRODUCTION* RELIGION. It is generally believed,	%ly INTRODUCTION* RELIGION. It is generally believed, upon the : %ly INTRODUCTION* RELIGION. It is generally believed,
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page24	606	book	nls			introduction: XT Pre/byterles, and Kirk Sejfions. ift, The General AJ	introduction: XT Pre/byterles, and Kirk Sejfions. ift, The General AJ
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page25	606	book	nls			*vi INTRODUCTION. Befides the eftablifned churches, there are a	*vi INTRODUCTION. Befides the eftablifned churches, there are a
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page26	606	book	nls			nor aboVe two hund- red	nor aboVe two hund- red
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page27	606	book	nls			INTRODUCTION. jcvi diſtriſt.	INTRODUCTION. jcvi diſtriſt.
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page28	606	book	nls			xvili INTRODUCTION. fophical works of Boethius is purely claffica	xvili INTRODUCTION. fophical works of Boethius is purely claffica
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page29	606	book	nls			and theodatinity of Buchaniari fe the moft claf	and theodatinity of Buchaniari fe the moft claf
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973437/	page	Page30	606	book	nls			INTRODUCTION. xi:c ion from his Majefty of 200l. per annum. Eve	INTRODUCTION. xi:c ion from his Majefty of 200l. per annum. Eve
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973438/	page	Page31	606	book	nls			xx INTRODUCTION. of their being a depifed people. The great Mr.	xx INTRODUCTION. of their being a depifed people. The great Mr.
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973438/	page	Page32	606	book	nls			xxi INTRODUCTION. *xi* CONSTITUTION. The ancient conſtituſion and	xxi INTRODUCTION. *xi* CONSTITUTION. The ancient conſtituſion and
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973438/	page	Page33	606	book	nls			I promife theſe three things to the Chriftian pe	I promife theſe three things to the Chriftian pe
gazetteer of Scotland	1803	1803	Dundee	/lustre/home/alt0/973438/	page	Page34	606	book	nls			wi: i n T D N L E Ti N model of the French nadirment to foun	wi: i n T D N L E Ti N model of the French nadirment to foun

Preprocessed Data

+ Text Mining queries



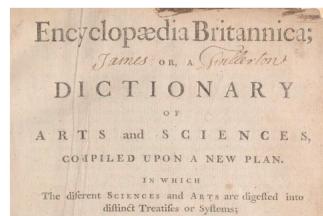
HDFS / ES/ PSQL

- [📄 georesolution_pages.py](#)
- [📄 keysearch_by_year.py](#)
- [📄 keysentence_concordance_by_yea...](#)
- [📄 normalize.py](#)
- [📄 window_concordance_by_date.py](#)
- [📄 keysearch_by_year.py](#)

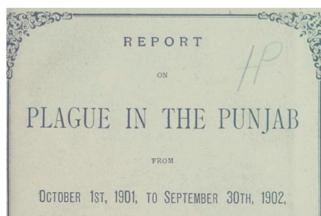
National Library Scottish – Uses Cases

Digitised collections

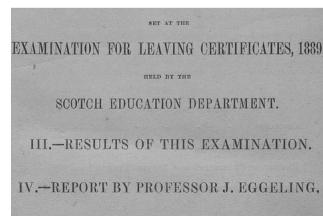
Download the ALTO, METS, image and plain text files for our digitised collections.



Encyclopædia Britannica, 1768-1860



A Medical History of British India

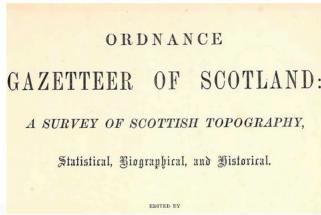


Scottish School Exam Papers, 1888-1963



The Spiritualist

<https://data.nls.uk/data/digitised-collections/a-medical-history-of-british-india/>



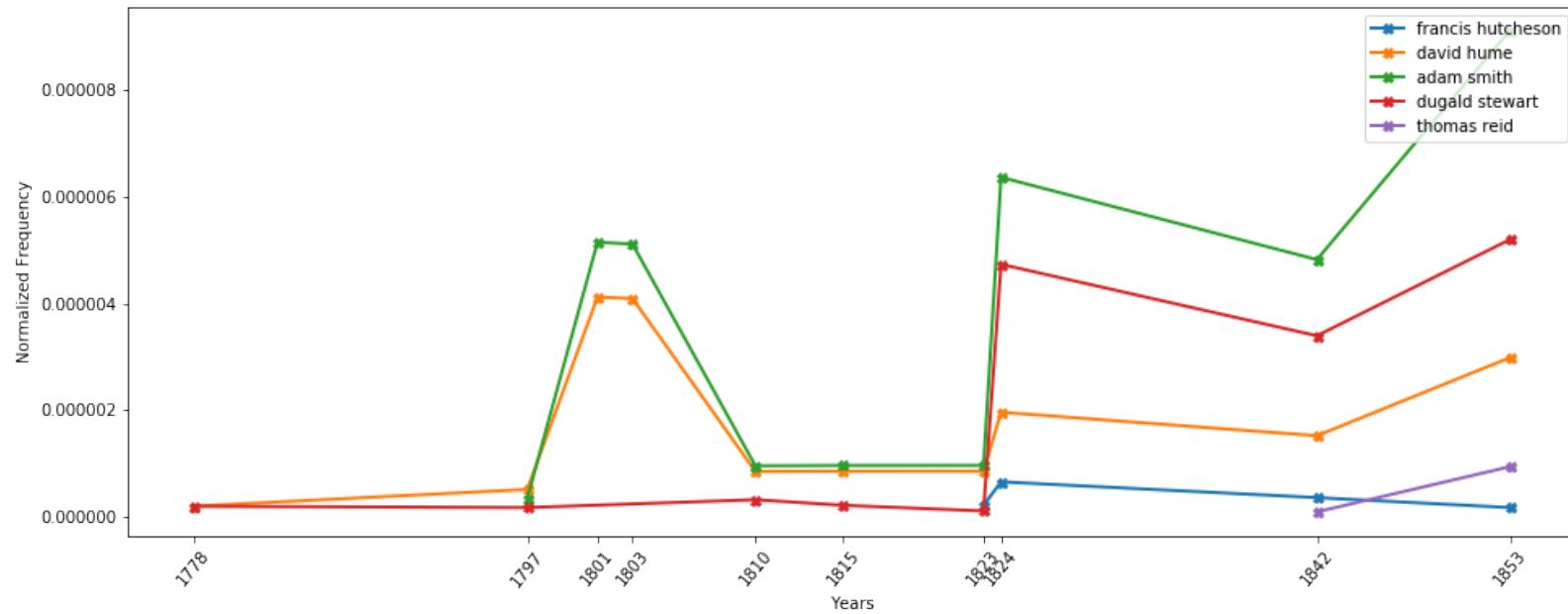
Gazetteers of Scotland



Ladies' Edinburgh Debating Society

National Library Scottish – Encyclopaedia Britannica

Exploring the popularity of Scottish Philosophers over time



N-gram for the normalised frequencies of the Scottish philosophers

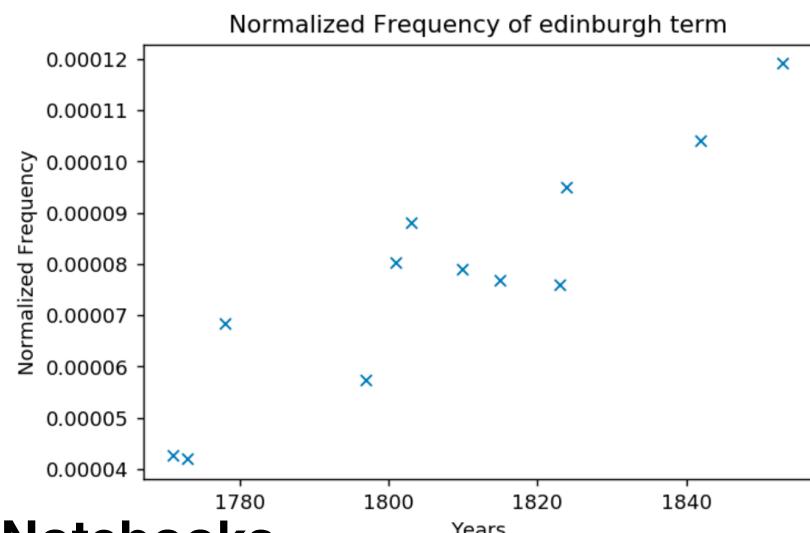
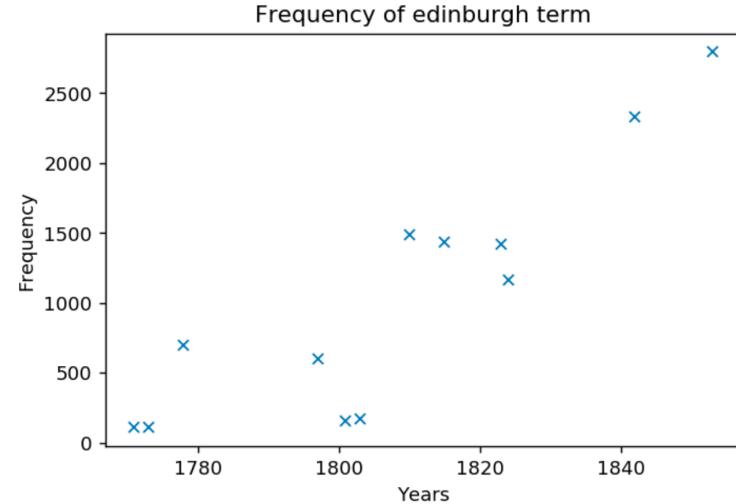
Check our Jupyter Notebooks

https://github.com/alan-turing-institute/defoe_visualization/tree/master/NLS

National Library Scottish – Encyclopaedia Britannica

Exploring the term “Edinburgh”

year	words_left	term	words_right
25	from', 'the', 'accounts', 'given', 'by', 'the', 'late', 'dr', 'whytt', 'of'	edinburgh	'and', 'others', 'it', 'should', 'appear', 'that', 'this', 'had', 'sometimes', 'happe...
34	stirling', 'and', 'frequently', 'the', 'three', 'principal', 'fortrefles', 'of', 'the'...	edinburgh	'stirling', 'and', 'dunbarton', 'the', 'last', 'heir', 'of', 'the', 'scottish', 'mona...
35	regent', 'mar', 'after', 'ihe', 'was', 'obliged', 'by', 'the', 'treaty', 'of'	edinburgh	'to', 'desist', 'from', 'wearing', 'the', 'arms', 'of', 'england', 'in', 'the']
94	considered', 'as', 'presages', 'of', 'a', 'pestilence', 'they', 'appear', 'annually', ...	edinburgh	'in', 'february', 'and', 'feedon', 'the', 'berries', 'of', 'the', 'mountainalh', 'they']
350	into', 'inches', 'the', 'flandard', 'is', 'kept', 'in', 'the', 'counleihambur', 'of'	edinburgh	'and', 'being', 'compared', 'with', 'the', 'engliah', 'yard', 'is', 'found', 'to']
437	with', 'in', 'duddinglilonloch', 'a', 'freshwater', 'lake', 'within', 'a', 'mile', 'of'	edinburgh	'in', 'france', 'it', 'stays', 'throughout', 'the', 'year', 'and', 'makes', 'a']
448	comprehensive', 'system', 'in', 'three', 'volumes', 'publiffed', 'by', 'mr', 'elliot'...	edinburgh	'upon', 'a', 'plan', 'approved', 'of', 'by', 'dr', 'monro', 'and', 'executed']
487	and', 'on', 'the', 'th', 'of', 'june', 'held', 'a', 'parliament', 'at'	edinburgh	'at', 'that', 'time', 'the', 'clergy', 'met', 'in', 'one', 'of', 'the']
488	mention', 'is', 'made', 'of', 'his', 'being', 'minister', 'of', 'pencaitland', 'near'	edinburgh	'in', 'but', 'we', 'find', 'nothing', 'said', 'there', 'or', 'any', 'where']
525	mild', 'difpolitions', 'and', 'peaceable', 'deportment', 'a', 'parliament', 'was', 'su...	edinburgh	'in', 'summer', 'and', 'the', 'duke', 'was', 'appointed', 'commiffioner', 'besides', ...



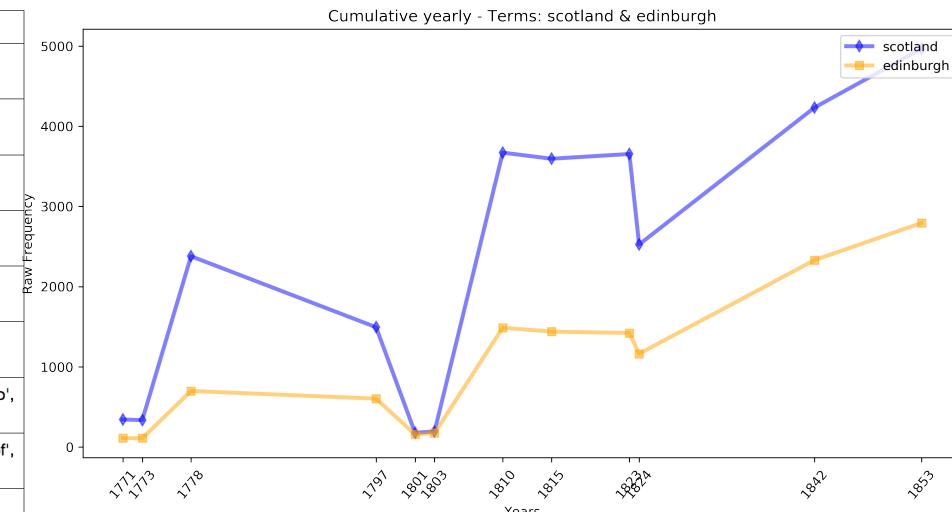
Check our Jupyter Notebooks

https://github.com/alan-turing-institute/defoe_visualization/tree/master/NLS

National Library Scottish – Encyclopaedia Britannica

Exploring the term “Edinburgh”

year	words_left	term	words_right
25	from', 'the', 'accounts', 'given', 'by', 'the', 'late', 'dr', 'whytt', 'of'	edinburgh	'and', 'others', 'it', 'should', 'appear', 'that', 'this', 'had', 'sometimes', 'happe...
34	stirling', 'and', 'frequently', 'the', 'three', 'principal', 'fortrefles', 'of', 'the'...	edinburgh	'stirling', 'and', 'dunbarton', 'the', 'last', 'heir', 'of', 'the', 'scottish', 'mona...
35	regent', 'mar', 'after', 'ihe', 'was', 'obliged', 'by', 'the', 'treaty', 'of'	edinburgh	'to', 'desist', 'from', 'wearing', 'the', 'arms', 'of', 'england', 'in', 'the']
94	considered', 'as', 'presages', 'of', 'a', 'pestilence', 'they', 'appear', 'annually', ...	edinburgh	'in', 'february', 'and', 'feedon', 'the', 'berries', 'of', 'the', 'mountainalh', 'they']
350	into', 'inches', 'the', 'andard', 'is', 'kept', 'in', 'the', 'counilehamber', 'of'	edinburgh	'and', 'being', 'compared', 'with', 'the', 'engliah', 'yard', 'is', 'found', 'to']
437	with', 'in', 'duddinglilonloch', 'a', 'freshwater', 'lake', 'within', 'a', 'mile', 'of'	edinburgh	'in', 'france', 'it', 'stays', 'throughout', 'the', 'year', 'and', 'makes', 'a']
448	comprehensive', 'system', 'in', 'three', 'volumes', 'publiffed', 'by', 'mr', 'elliot'...	edinburgh	'upon', 'a', 'plan', 'approved', 'of', 'by', 'dr', 'monro', 'and', 'executed']
487	and', 'on', 'the', 'th', 'of', 'june', 'held', 'a', 'parliament', 'at'	edinburgh	'at', 'that', 'time', 'the', 'clergy', 'met', 'in', 'one', 'of', 'the']
488	mention', 'is', 'made', 'of', 'his', 'being', 'minister', 'of', 'pencaitland', 'near'	edinburgh	'in', 'but', 'we', 'find', 'nothing', 'said', 'there', 'or', 'any', 'where']
525	mild', 'difpolitions', 'and', 'peaceable', 'deportment', 'a', 'parliament', 'was', 'su...	edinburgh	'in', 'summer', 'and', 'the', 'duke', 'was', 'appointed', 'commiffioner', 'besides', ...



Check our Jupyter Notebooks

https://github.com/alan-turing-institute/defoe_visualization/tree/master/NLS

National Library Scottish – Encyclopaedia Britannica

[Encyclopaedia Britannica: or, A dictionary of arts and sciences]:
archive_name: /home/tmd/datasets/eb_test/144850366
articles:
4 ✓ **ACQUEST**: or Acquist, in law, signifies goods got by purchase or donation. See CONQUEST.
5 ✓ **ACQUI:** "a town of Italy, in the Duchy of Montferrat, with a bishop[u2019]s see, \
6 ✓ and commodious baths. It was taken by the Spaniards in 1745, and retaken by \
7 ✓ the Piedmontese in 1746; but after this, it was taken again and disfranchised \
8 ✓ by the French, who afterwards forsook it. It is seated on the river Bormio, \
9 ✓ 25 miles[N. W.] Genoa, and 30S. Cafal, 8. 30.E. 44. 50.
10 ✓ . lat." 10.
11 ✓ **ACQUIESCENCE**: in commerce, is the consent that a person gives to the determination
12 given either by arbitration, orbycausal
13 ✓ **ACQUITETANDIS plegis**; in the English law, is a writ that lies for a surety, against
14 a creditor, who refuses to acquit the complainant after the debt is paid.
15 ✓ **ACQUITTAINTA de Jheris et hundredis**: in England, signifies the privilege of being
16 free from suit and service in firlies and hundreds.
17 ✓ **ACQUISITION**: in general, denotes the obtaining or procuring something. Among lawyers,
18 it is used for the right or title to an estate got by purchase or dona-
19 **ACQUITARE**: in ancient law-books, signifies to discharge or pay off the debts
20 of a person deceased.
21 ✓ **ACQUITTALE**: a discharge, deliverance, or setting of a person free from the guilt
22 or suspicion of an offence.
23 **ACQUITTANCE**: a release or discharge in writing for a sum of money.
24 ✓ **ACRA**: a town of Africa, on the coast of Guinea, where the English, Dutch, and
25 Danes, have strong sorts, and each sort its particular village, o. 2.W. 5. o.
lat.
26 ✓ **ACRASIA**: among physicians, signifies the predominancy of one quality over another.
27 **ACRE**: "or Acra, a sea-port town in Syria. It was formerly called Ptolemais, and \
28 ✓ \ is a bishop[u2019]s see. It was very famous in the time of the crusades, and \
29 ✓ underwent several sieges both by the Chrif Hans and Saracens. It is now an \
30 ✓ inconconsiderable town, being entirely supported by its harbour, which, is frequented \
31 ✓ by ships of several nations. It is 20 miles[S. Tyre, and 37N. Jerusalem, 39. \
32 ✓ 25.E. 32. 40. lat. Acre, in the Mogul[u2019]s dominions, the same with lack, \
33 ✓ and signifies the sum of 100,000 rupees; the rupee is not of the value of the
34 ✓ French crown of 3 livres, or 30 florins of Holland; the 100 lacks of rupees make \
35 ✓ a couron in Indofaan, or 10,000,000 rupees; the pound Sterling is about 8
36 ✓ rupees; according to which proportion, a lack of rupees amounts to 12,500
37 ✓ pounds Sterling. Acre, a measure of land used in several provinces of France, \
38 ✓ particularly in Normandy. It is larger or less according to the different \
39 ✓ places ; but commonly contains 160 perches. The Acre of woods in France, consists \
40 ✓ of four rods, called verg es; the rod is 40 perches, the perch 24 feet, \
41 ✓ the foot 12 inches, the inch 12 lines. Acre, the universal measure of land, \
42 ✓ in Britain. An acre in England contains 4 square rods, a rod 40ACR or poles \
43 ✓ of 16X32 feet each byttafe. Yet this measure does not prevail in all parts \
44 ✓ of England, as the length of the pole varies in different counties, and is \
45 ✓ called cuijosity measure, the difference running from the 16feet to 28. The \
46 ✓ acre is also divided into xo square chains, of 22 yards each, that is 4840 \
47 ✓ square yards. An acre in Scotland contains 4 square rods ; 1 square rood \
48 ✓ is 40 square falls; 1 square fall, 36 square ell; 1 square ell, 9 square \
49 ✓ feet, and 73 square inches; 1 square foot, 144 square inches. The Scots acre \
50 ✓ is also divided into 10 square chains ; the measuring chrfn in Iboe are 24 \
51 ✓ ell in length, divided into 100 links, each link 8 r [u2019]c 71 and so1 \
52 ✓ chain will contain 10,000 square links. The English (latute acre is about \
53 ✓ 3 rods) and 6 falls standard measure of Scotland."
54 ✓ **ACREME**: in old law-books, signifies ten acres of land.
55 ✓ **ACRIBEA**: signifies great accuracy.
56 ✓ **ACRIDID**: a name for any thing that is of a (harp or pungent taste.
57 ✓ **ACRIDOPHAGI**: "signifies loculi-eaters. It has been much disputed whether the inhabitants \
58 ✓ of Arabia, Ethiopia, <bc ever eat locusts. We shall give the substance of \
59 ✓ what Haffequill says on this subject; who travelled in Syria and Egypt so \
60 ✓ late as the year 1752. This ingenious gentleman, who travelled with a view \
61 ✓ to improve natural history, informs us, that he asked Franks, and many other \
62 ✓ people who had lived long in these countries, whether they had ever heard \
63 ✓ that the inhabitants of Arabia and Ethiopia, &c. used locusts as food. They \
64 ✓ answered that they had. He likewise asked the same question of Armenians, \
65 ✓ Coptes, and Syrians, who lived in Arabia, and had travelled in Syria and \
66 ✓ near the Red-sea; some of whom said they heard of such a practice, and others \
67 ✓ that they, had often seen the people eat these infeds. In the last obtained \
68 ✓ complete satifaction on this head from a learned Heck at Cairo, who had \
69 ✓ lived six years in Mecca. This gentleman told him, in presence of M. Grand, \
70 ✓ the principal French interpreter at Cairo, and others, that a famine frequently \
71 ✓ rages at Mecca when there is a scarcity of corn in Egypt, which obliges the \
72 ✓ inhabitants to live upon coarser food than ordinary : That when corn is scarce, \
73 ✓ the Arabians grind the loculls in hand-mills, or Hone mortars, and bake them \
74 ✓ into cakes, and use these cakes in place of bread: That he has frequently \
75 ✓ seen locusts used by the Arabians, even when there was no scarcity of corn; \
76 ✓ but then they boil them, itew them with butter, and make them into a kind \
77 ✓ of fricoffee, which he says is not disagreeably tailed; for he had sometimes \
78 ✓ tailed these locutt-friccoffees out of curiosity. From this account, we may \
79 ✓ see the folly of that dispute among divines about the nature of St John[u2019]s \
80 ✓ food in the wilderness. Some of them say that loculls were the fruits of certain \
81 ✓ trees, others that they were a kind of birds, <bc.; but those who adhered \
82 ✓ to the former opinion, were not able to give any reason for their opinion."

A C R (20) A C R

ACQUEST, or Acquist, in law, signifies goods got by purchase or donation. See CONQUEST.

ACQUI, a town of Italy, in the Duchy of Montferrat, with a bishop's see, and commodious baths. It was taken by the Spaniards in 1746; but after this, it was taken again and dismantled by the French, who afterwards forsook it. It is seated on the river Bormio, 25 miles N. W. of Genoa, and 30 S. of Cafal, 8. 30. E. long. 44. 0. lat.

ACQUIESCENCE, in commerce, is the consent that a person gives to the determination given either by arbitration, or by a court.

ACQUIETANDIS *plegis*, in the English law, is a writ that lies for a fury, against a creditor, who refuses to acquire the complainant after the debt is paid.

ACQUIETANTIA de *sistit et hundredis*, in England, signifies the privilege of being free from suit and service in shires and hundreds.

ACQUISITION, in general, denotes the obtaining or procuring something. Among lawyers, it is used for the right or title to an estate got by purchase or donation.

ACQUITARE, in ancient law-books, signifies to discharge or pay off the debts of a person deceased.

ACQUITTAL, a discharge, deliverance, or setting of a person free from the guilt or suspicion of an offence.

ACQUITTANCE, a release or discharge in writing for a sum of money.

ACRA, a town of Africa, on the coast of Guinea, where the English, Dutch, and Danes, have strong forts, and each fort its particular village, o. 2. W. long. 5. 0. lat.

ACRASIA, among physicians, signifies the predominancy of one quality over another.

ACRE, or ACRA, a sea-port town in Syria. It was formerly called *Ptolemais*, and is a bishop's see. It was very famous in the time of the crusades, and underwent several sieges both by the Christians and Saracens. It is now an inconsiderable town, being entirely supported by its harbour, which is frequented by ships of several nations. It is 20 miles S. of Tyre, and 37 N. of Jerusalem, 39. 25. E. long. 32. 40. lat.

ACRE, in the Mogul's dominions, the same with lack, and signifies the sum of 100,000 rupees; the rupee is of the value of the French crown of 3 livres, or 30 sols of Holland; an 100 lacks of rupees make a couron in Indostan, or 10,000,000 rupees; the pound Sterling is about 8 rupees; according to which proportion, a lack of rupees amounts to 12,500 pounds Sterling.

ACRE, a measure of land used in several provinces of France, particularly in Normandy. It is larger or less according to the different places; but commonly contains 160 perches.

The ACRE of woods in France, consists of four rods, called *verg e*; the rod is 40 perches, the perch 24 feet, the foot 12 inches, the inch 12 lines.

ACRE, the universal measure of land in Britain. An acre in England contains 4 square rods, a rod 40

perches or poles of 16¹/₂ feet each by statute. Yet this measure does not prevail in all parts of England, as the length of the pole varies in different counties, and is called *cumfermey measure*, the difference running from the 16¹/₂ feet to 28. The acre is also divided into 10 square chains, of 22 yards each, that is 480 square yards. An acre in Scotland contains 4 square rods; 1 square rod is 40 square falls; 1 square fall, 36 square ells; 1 square ell, 9 square feet, and 73 square inches; 1 square foot, 144 square inches. The Scots acre is also divided into 10 square chains; the measuring chain should be 24 ells in length, divided into 100 links, each link 8¹/₂ inches; and so 1 square chain will contain 10,000 square links.

The English statute acre is about 3 rods and 6 falls standard measure of Scotland.

ACREME, in old law-books, signifies ten acres of land.

ACRIBEIA, signifies great accuracy.

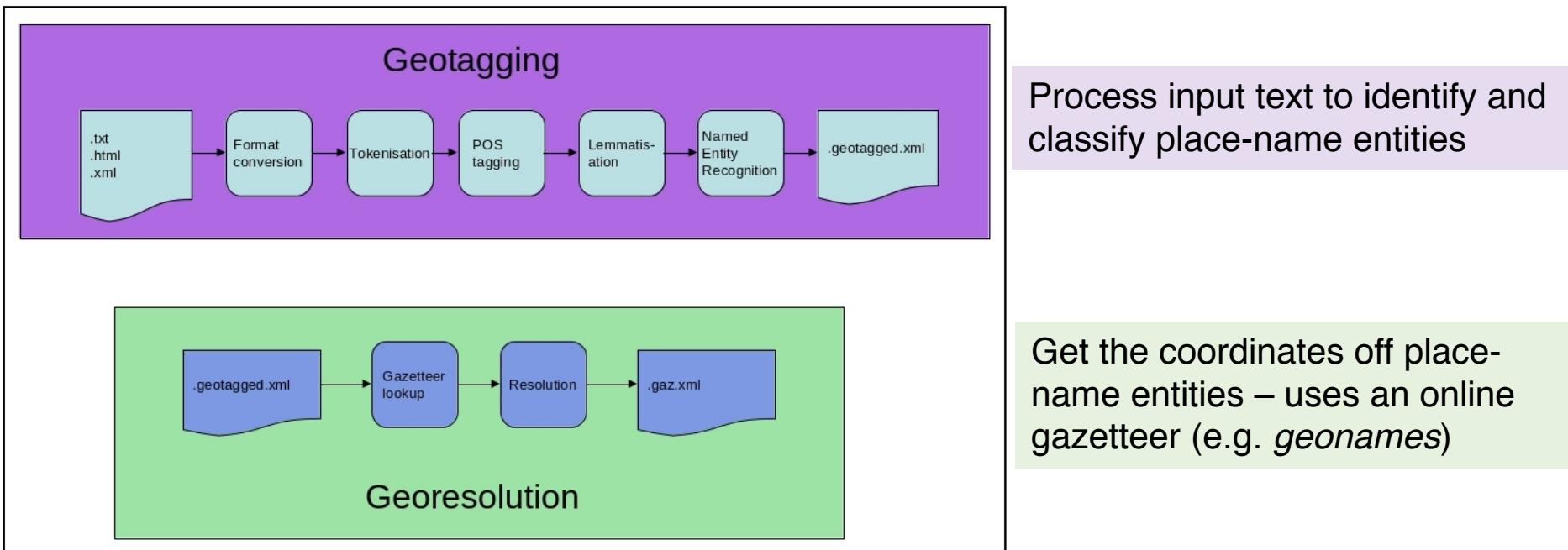
ACRID, a name for any thing that is of a sharp or pungent taste.

ACRIDOPHAGI, signifies locust-eaters. It has been much disputed whether the inhabitants of Arabia, Ethiopia, &c. ever eat locusts. We shall give the substance of what Hasselquist says on this subject, who travelled in Syria and Egypt to late as the year 1752. This ingenious gentleman, who travelled with a view to improve natural history, informs us, that he asked Franks, and many other people who had lived long in these countries, whether they had ever heard that the inhabitants of Arabia and Ethiopia, &c. used locusts as food. They answered that they had. He likewise asked the same question of Armenians, Coptes, and Syrians, who lived in Arabia, and had travelled in Syria and near the Red-sea; some of whom said they heard of such a practice, and others that they had often seen the people eat these insects. He at last obtained complete satisfaction on this head from a learned sleekt at Cairo, who had lived five years in Mecca. This gentleman told him, in presence of M. le Grand, the principal French interpreter at Cairo, and others, that a famine frequently rages at Mecca when there is a scarcity of corn in Egypt, which obliges the inhabitants to live upon coarser food than ordinary: That when corn is scarce, the Arabians grind the locusts in hand-mills, or stone mortars, and bake them into cakes, and use these cakes in place of bread: That he has frequently seen locusts used by the Arabians, even when there was no scarcity of corn; but then they boil them, stew them with butter, and make them into a kind of friezafee, which he says is not disagreeably tasted; for he had sometimes tasted their locust-sarcasses out of curiosity. From this account, we may see the folly of that dispute among divines about the nature of St John's food in the wilderness. Some of them say that locusts were the fruits of certain trees, others that they were a kind of birds, &c.; but those who adhered to the literal meaning of the text were at least the most orthodox, although their arguments were perhaps not so strong as they might have been, had they had an opportunity of quoting such an author as Hasselquist,

National Library Scottish – Scottish Gazetteers

Geoparsing the historical Gazetteers of Scotland

Edinburgh Geoparser + Defoe



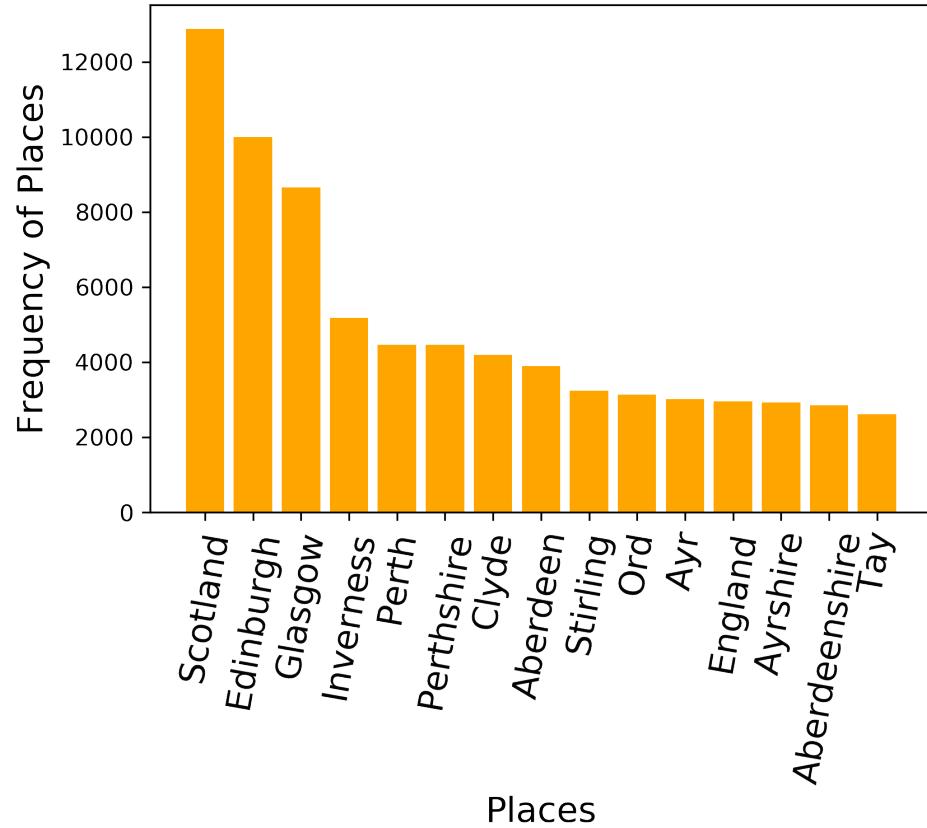
Edinburgh Geoparser Pipeline

Collaboration with the **Language Technology Group** – University of Edinburgh

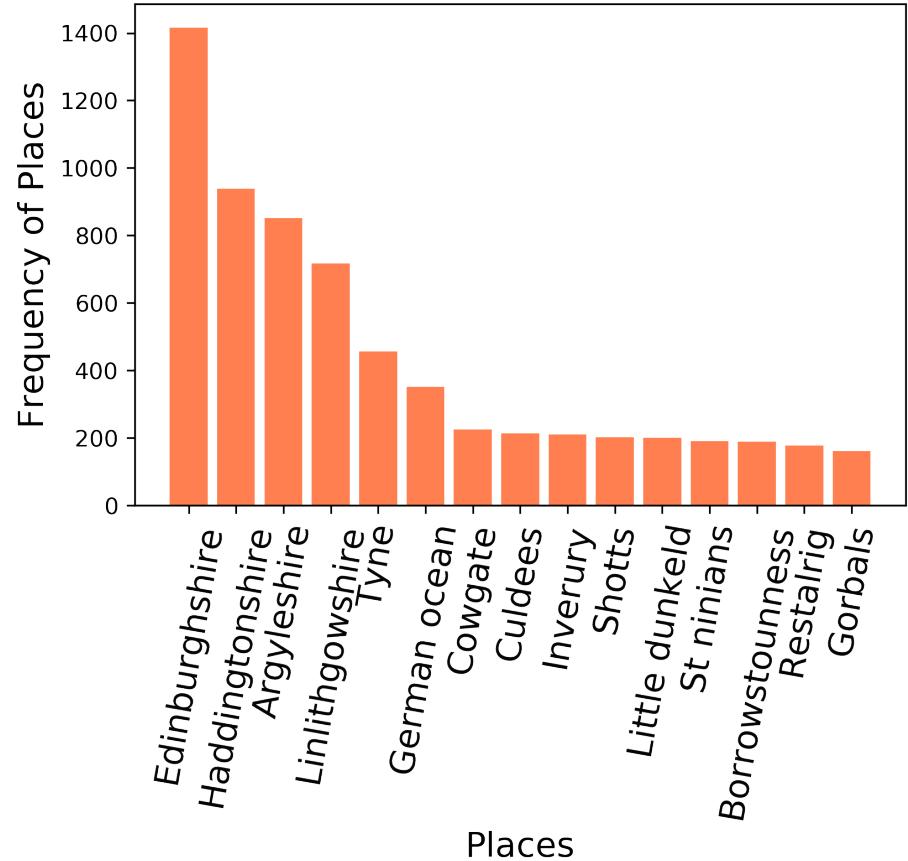
<https://www.ltg.ed.ac.uk/software/geoparser/>

National Library Scottish – Scottish Gazetteers

15 Places most mentioned using
the Original Geoparser across all Scottish Gazetteers



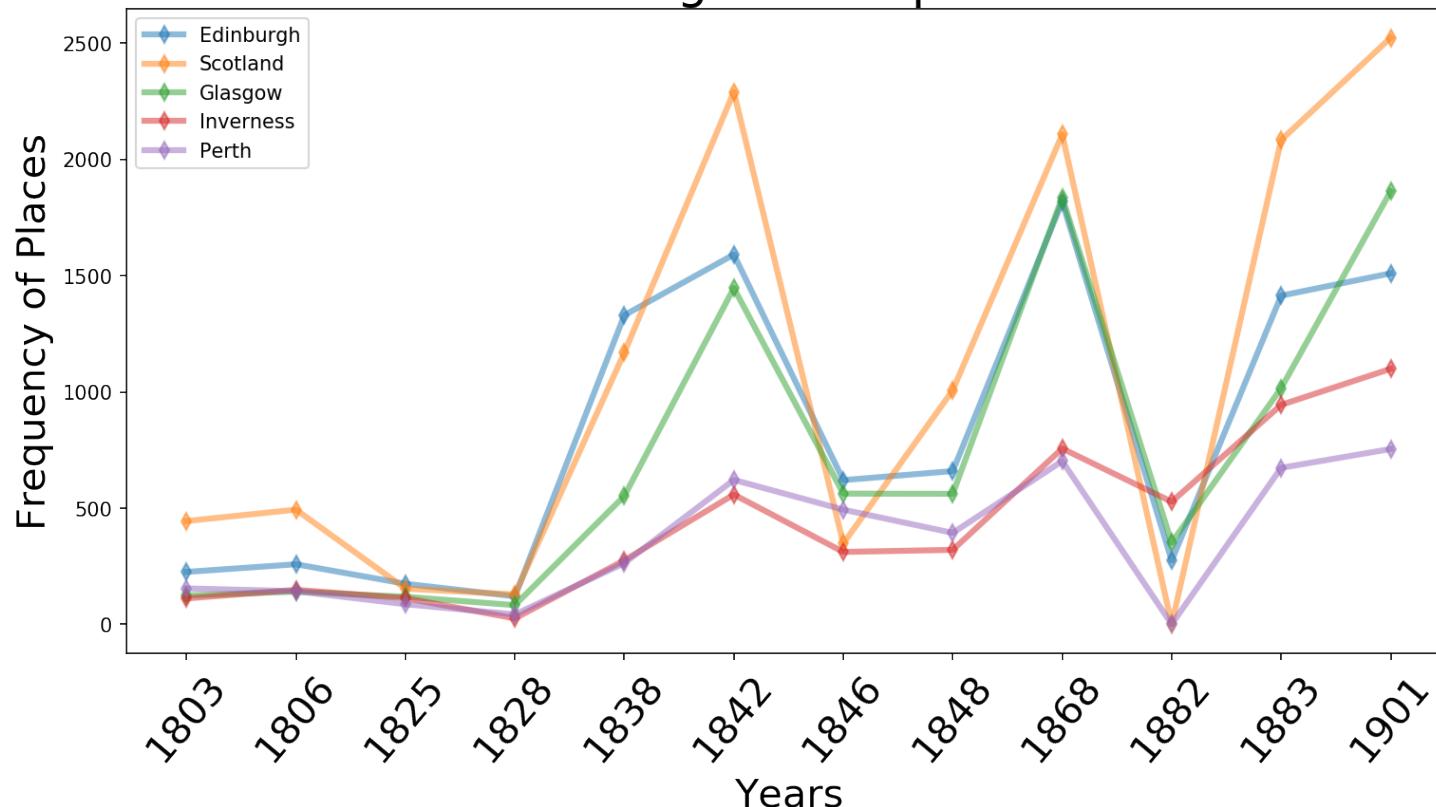
15 Places most mentioned but not resolved using
the Original Geoparser across all Scottish Gazetteers



Check our Jupyter Notebooks

National Library Scottish – Scottish Gazetteers

Yearly cumulative frequencies using
the Original Geoparser



Publication: Geoparsing the Historical Gazetteers of Scotland: Accurately Computing Location in Mass Digitised Texts, 8th Workshop on the Challenges in the Management of Large Corpora, 16th May 2020

DATA, CULTURE & SOCIETY



[Home](#) / [CDCS Text Mining Lab - Call for Projects](#)

CDCS TEXT MINING LAB - CALL FOR PROJECTS

The Centre for Data, Culture & Society is looking for humanities and social science researchers who can ask complex questions of large-scale data sets.

Deadline for expressions of interest: April 20th - Applications are now closed.

This call is open to research active staff and PhD students within CAHSS at the University of Edinburgh.

****Please note that we envisage being able to conduct this activity remotely.****





DATA, CULTURE & SOCIETY



[Home](#) / [CDCS Text Mining Lab - Call for Projects](#)

CDCS TEXT MINING LAB - CALL FOR PROJECTS

DATA SETS

- BRITISH LIBRARY BOOKS

Over 68,000 books from the 16th to the 19th century, covering geography, philosophy, history, poetry and literature in a variety of languages.

- BRITISH LIBRARY NEWSPAPERS

1TB of digitised British newspapers from the 18th to the early 20th Century

- TIMES DIGITAL ARCHIVE

All the articles in 69,699 volumes of The Times newspaper between 1785 and 2009.

- PAPERS PAST: NEW ZEALAND AND PACIFIC NEWSPAPERS

Over 5 million pages of New Zealand and Pacific newspapers from the 19th and 20th Centuries.

- GAZETTEERS OF SCOTLAND

20 volumes of the most popular 19th Century gazetteers of Scotland, including detailed historical and geographic information about each place.

- ENCYCLOPAEDIA BRITANNICA 1768 - 1860

The first 8 volumes of the Encyclopaedia Britannica, issued from 1768-1860, comprising a total of 143 volumes. 155,388 pages, 166m words.

Coming soon:

- JISC MEDICAL HERITAGE LIBRARY
- HANSARD ARCHIVE
- STATISTICAL ACCOUNTS OF SCOTLAND
- DIGITISED THESES FROM EDINBURGH UNIVERSITY LIBRARY

If we don't have what you're looking for in term of data sets, you can propose a new data set for use with *defoe*.

Conclusions

- New digital toolbox for mining historical data.
- Text analyses across large collections in parallel.
- Rich set of text mining queries.
- NLP prepossessing techniques to mitigate OCR errors.
- Portability on different computing environments
- Compatibility with different storage systems
- Integration with other NLP tools – e.g. Edinburgh Geoparser

"All this work provides the means to search across large scale datasets and to return results for further analysis and interpretation by historians."



Analysing Digital Historical Textual Data at Scale with Defoe Toolbox

Dr. Rosa Filgueira, EPCC,
University of Edinburgh
Email: rosa.filgueira@ed.ac.uk