

# Data Mining for Historians

Professor Bob Shoemaker  
University of Sheffield  
12 May 2010

*“instead of looking for a needle in the haystack,  
an effective text mining tool will try and show  
you the shape of the haystack and tell you the  
words you might want to find”*  
—Geoffrey Rockwell



## Home Page

## Search

## About the Proceedings

## Historical Background

## The Project

## Copyright & Citation Guide

## Contact

## ON THIS DAY IN... 1736

William Rine allegedly pulled a man off of his horse, robbed him, stripped him naked, and tied him to a tree. [read more](#)

## The Proceedings of the Old Bailey, 1674-1913

A fully searchable edition of the largest body of texts detailing the lives of non-elite people ever published, containing 197,745 criminal trials held at London's central criminal court.

To search the Proceedings use the boxes on the right or go to the [Search Pages](#).

## What's New?

Trials between November 1834 and April 1913 and the **Ordinary of Newgate's Accounts** between 1690 and 1772 have now been added. See [What's New](#).

## Ordinary's Accounts

The website now includes the texts of all [Ordinary of Newgate's Accounts](#) published between 1690 and 1772. These richly detailed narratives of the lives and deaths of convicts executed at Tyburn have been linked to the relevant

## SEARCH

*the Proceedings*

Keyword(s)

Reference No.

Search In

<All Text>



SEARCH

[More Search Options](#)

CENTRAL CRIMINAL COURT.

SESSIONS PAPER.

*In this Section...*

## Search Home

Personal Details

Ordinary's Accounts,  
1676-1772

Proceedings by date

Ordinary's Accounts by date

Statistics

Custom Search

Associated Records, 1674-1834

Place and Map Search,  
1674-1834

### Security Specialists

Get A home security  
system In Grimsby.

[www.abatis-fs.co.uk](http://www.abatis-fs.co.uk)

Ads by Google

### Home & Family Protection

## Search Home

The boxes below allow you to search the whole of the **Proceedings** and all published **Ordinary Accounts** (for the period 1679 to 1772). You may combine keyword searches with queries on the information including **surname**, **crime**, and **punishment**. The default setting allows you to search the full text of all the documents available on this website. This page should be used for basic and simple searches. Please refer to the other pages listed to your left for more search options.

**Keyword(s)**



**Surname**



**Given Name**



**Alias**



**Offence**



**Verdict**



**Punishment**



**Search In**



**Time Period**

*From (month/year)* *To (month/year)*

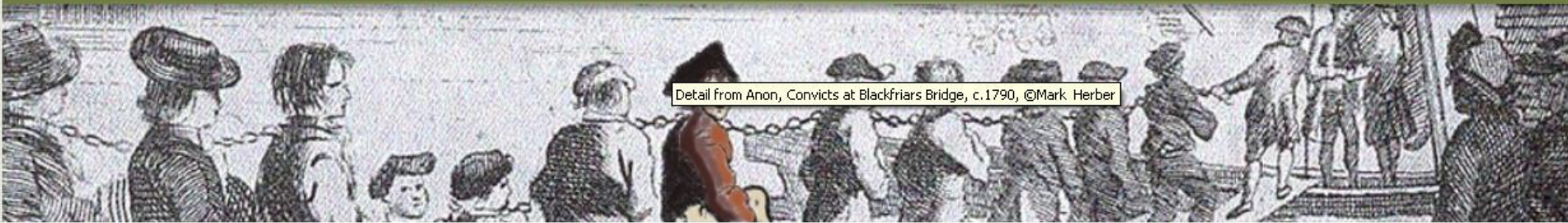


**Reference Number**



SEARCH





### [Home](#)

[Search](#)[Login / Register](#)[Background](#)[Lives](#)[Project](#)[Legal](#)[Wiki](#)[Contact](#)

## Sources for London Lives

A fully searchable edition of 240,000 manuscripts from eight archives and twelve datasets, giving access to over 3.4 million names.

**To search London Lives use the boxes on the right or go to the Search Pages.**

## Help us Construct London Lives

You can use this site to create biographies of eighteenth-century Londoners. The [person search page](#) allows you to assemble a set of documents related to a single person, and the [project wiki](#) allows you to write and share their biography. You will need to register and login in order to make use of these features. See the Featured Life on this page, and the [Lives section](#).

## London Life in the Eighteenth Century

What was it like to live in the world's first million person city? Crime, poverty, and illness; apprenticeship, work, politics and money; how people voted, lived and died; all this and more can be found in these documents. For more

## SEARCH

Keyword

Surname

Given Name

Reference ID

From

1680

To

1820

**SEARCH**[More Search Options](#)

## Advertisements

[Google code goes here]

## Featured Life

**Sophie Pringle, Executed 1787**

## Browse by Document Type

**London Lives** includes 39 individual **Document Types** from fourteen archives, and 15 **Additional Datasets**. By clicking on an individual **Archive**, you will be taken to a page which lists all the document types from that archive. By then clicking on a decade for which those documents are available, you will be taking to a listing of that document type year by year.

### Parish Archives

- [St Botolph Aldgate](#)
- [St Clement Danes](#)
- [St Dionis Backchurch](#)

### Criminal Records

- [Bridewell Royal Hospital](#)
- [Home Office](#)
- [Old Bailey Sessions](#)
- [Old Bailey Proceedings](#)
- [Ordinary's Accounts](#)
- [City of London Sessions](#)
- [Middlesex Sessions](#)
- [Westminster Sessions](#)

### Coroners' Records

- [City of London Coroners](#)
- [Middlesex Coroners](#)
- [City of Westminster Coroners](#)

### Hospital and Guild Records

- [Carpenters' Company](#)
- [St Thomas's Hospital](#)

[illegible]





# **Connected Histories: Sources for Building British History, 1500-1900**

- **Funded by JISC (UK) e-Content Capital Programme**
- **Partnership between:**
  - **University of Sheffield (Humanities Research Institute)**
  - **University of Hertfordshire**
  - **Institute of Historical Research, London**





In this Section...

[Home Page](#)

[Search](#)

[About the Proceedings](#)

[Historical Background](#)

### The Proceedings of the Old Bailey, 1674-1913

A fully searchable edition of the largest body of texts detailing the lives of non-elite people ever published, containing 197,745 criminal trials held at London's central criminal court.

SEARCH

the Proceedings

Keyword(s)

Reference No.

## HOUSE OF COMMONS PARLIAMENTARY PAPERS



■ [SEARCH](#) ■ [BROWSE](#) ■ [INFORMATION RESOURCES](#) ■ [MY ARCHIVE](#) ■ [HELP](#)



Welcome to the House of Commons Parliamentary Papers (HCPP).

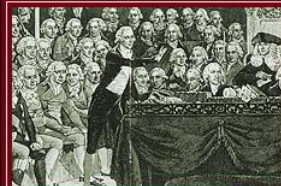
With all collections enabled, HCPP now includes over 200,000 House of Commons sessional papers from 1715 to the present, with supplementary material back to 1688. HCPP delivers page images and searchable full text for each paper, along with detailed indexing. Find out more [about House of Commons Parliamentary Papers](#).

You have access to the following collections:

18th century (1688-1834)

19th century (1801-1900)

20th century (1901-2003/04 session)



ARCTIC EXPEDITION UNDER SIR JOHN FRANKLIN. 127

No. 11 (D.)

Copy of a LETTER from Dr. McCormick to the Secretary of the Admiralty.

11, Apsley Cottages, Twickenham Green,  
20 February 1850.

Sir,  
I have been to transmit herewith, for the approval of my Lords Commissioners

No. 11.  
Dr. McCormick's  
Plan for a Boat  
Expedition.

originsnetwork

The premier resource for tracing your British and Irish ancestors

my Origins [LOGIN](#) [?](#)

[SUBSCRIPTION OPTIONS](#) [HELP & RESOURCES](#) [ORIGINS SHOP](#) [DISCUSSION](#) [WHAT'S NEW](#)

New users [SUBSCRIBE](#) [?](#)



Welcome to Origins Network

► Search the **richest source** of British and Irish genealogy online

► View **comprehensive record collections**, dating back to the 13th century, most not available anywhere else on the internet

► Browse **rare and vintage** photos, maps and books

► Plus over 3000 books and cds available in the **Origins Shop**

[FIND OUT MORE](#) [?](#)

[SUBSCRIBE](#) [?](#)

TRY A FREE SEARCH BY NAME BY PLACE

Begin your family history research below

Last name

NameX [Close variants](#) [?](#)

First name

NameX [Close variants](#) [?](#)

[SEARCH](#) [?](#)

[SUBSCRIPTION OPTIONS](#)

[GIFTS / VOUCHERS](#)

[LIBRARY LICENCE](#)

[SOG MEMBERS](#) [?](#)



#### Total Access subscription

Best value membership!  
Search **both** British and Irish collections  
and Libraries for one low price.  
[Find out more](#)



#### British Origins subscription

Looking for ancestors in Britain? Seek  
them out in Origins richest source of  
British genealogy online.  
[Find out more](#)

WHAT'S NEW September 2006 [JOIN MAILING LIST](#)

#### NEW! Dorset Marriage Index 1538-1856

Over 150,000 marriages transcribed by Somerset &  
Dorset Family History Society. Available to **British  
Origins** and **Total Access** subscribers.  
[Search now](#)

[Additional Directories of Ireland](#)

## John Styrpe's A SURVEY OF THE CITIES OF London and Westminster

[home](#) | [transcriptions](#) | [parishes](#) | [illustrations](#) | [search](#) | [help](#) | [about](#) 

An electronic edition of John Styrpe's

## A SURVEY OF THE CITIES OF London and Westminster

Version 1.0 (ISBN: 0-9542608-9-9)



John Stow's Elizabethan classic, *A Survey of  
London*, was first published in 1598, with a  
second edition following in 1603. Stow (c.

See also:

- Old Bailey Online and London Lives
- British History Online
- Burney Newspaper Collection
- Origins Network (genealogical database)
- Parliamentary Papers
- Clergy of the Church of England Database
- Charles Booth Online Archive
- Collage
- John Johnson Collection of Printed Ephemera (provisional)
- EEBO and ECCO (provisional)
- *Others subject to negotiation (more can be easily added)*

# ***Connected Histories Methods***

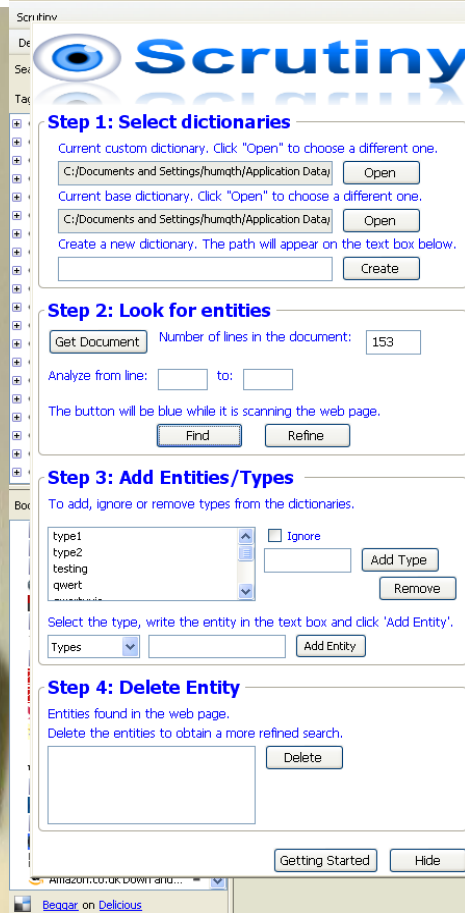
- **A federated search facility, sitting as umbrella over the distributed sites**
- **Search by keyword, name, place, date**
- **Use existing metadata and tagged information where available**
- **Where information not tagged, data will be obtained by web crawling and natural language processing will be used to identify key entities**
- **Data saved as RDF files held centrally**

# ***Connected Histories Outcomes***

- Search engine on Connected Histories website ([www.connectedhistories.org](http://www.connectedhistories.org))
- Search results point users to original sites (with login for commercial sites)
- Background information
- User engagement
- Website launch early 2011



Scrutiny ; fu  
not be imagined



[Home](#) | [Search](#) | [About The Proceedings](#) | [Historical Background](#) | [The Project](#) | [User Wiki](#) | [Contact](#)



**Available from: <https://addons.mozilla.org/en-US/firefox/addon/75085>**

# *Data Mining with Criminal Intent*

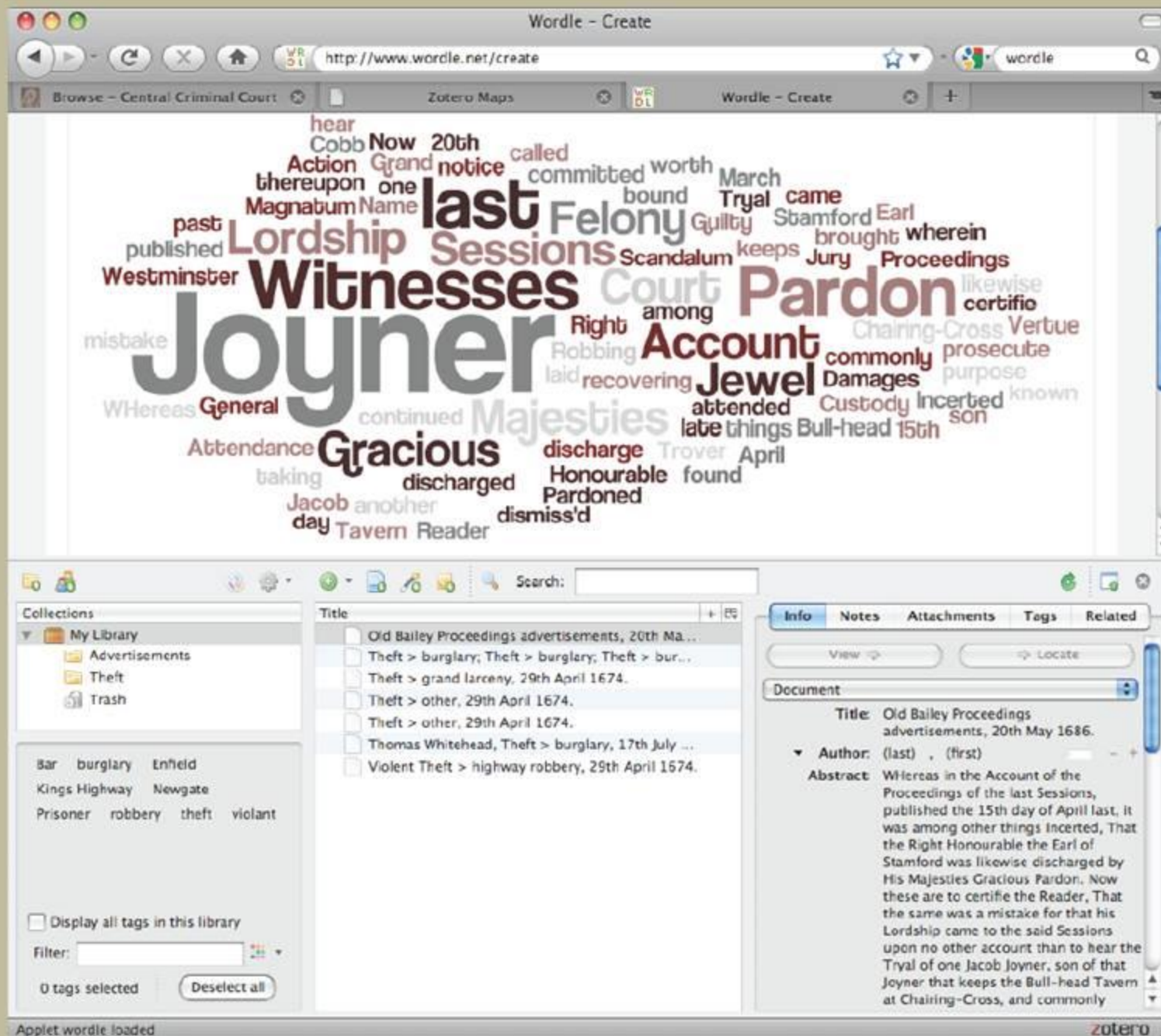
- Project Partners
  - Tim Hitchcock (University of Hertfordshire)
  - Robert Shoemaker (University of Sheffield)
  - Michael Pidd (Humanities Research Institute, University of Sheffield)
  - Daniel Cohen (Center for History and New Media, George Mason University)
  - Geoffrey Rockwood (University of Alberta)
  - William Turkel (University of Western Ontario)

# *Data Mining with Criminal Intent*

- Tools:
  - Zotero (Firefox Plug-in)
  - TAPoR toolkit (including Voyeur)
  - Compression Analysis



## Word Cloud Visualization of Old Bailey text via Zotero



# Invoking the TAPoR Frequency Tool with the Old Bailey via Zotero

**Summary:** There are 390 unique words other than those in the stop list, there are 928 words other than those in the stop list. There are 2037 words in total including the stop words.

Words	Distribution	Counts
persons		17
cures		14
england		9
diseases		9
teeth		9
water		8
city		8
mouth		8

**Zotero Collections:** My Library, Advertisements, Theft, Trash

**Tags:** Bar, burglary, Enfield, Kings Highway, Newgate, Prisoner, robbery, theft, violent

☐ Display all tags in this library

Filter:

0 tags selected

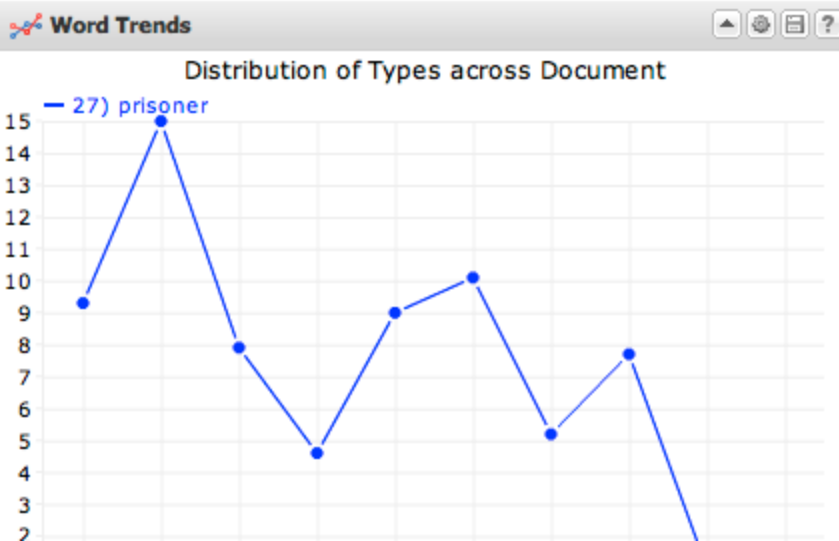
**Title**

- ☐ Old Bailey Proceedings advertisements, 20th Ma...
- ☐ Theft > burglary; Theft > burglary; Theft > bur...
- ☐ Theft > grand larceny, 29th April 1674.
- ☐ Theft > other, 29th April 1674.
- ☐ Theft > other, 29th April 1674.
- ☐ Thomas Whitehead, Theft > burglary, 17th July ...
- ☐ Violent Theft > highway robbery, 29th April 1674.

No items selected

Done zotero

Words in the Entire Corpus		
Word	Count	Trend
said	20,303	
prisoner	17,443	
mr	10,032	
went	8,499	
saw	7,961	
street	6,543	
came	5,887	
told	5,785	
know	5,769	
road	5,606	
man	5,275	
asked	5,195	
money	4,766	
got	4,647	



Summary Corpus

- There are 41 documents in this corpus with a total of **2,535,564 words** and **35,475 unique words**.
- Longest documents** (by words): T19101011CLR\_SUP\_DONE (111,217), T19120227NW\_SUP\_DONE (103,304). Shortest documents: T19120722GS\_SUP\_DONE (22,965), T19110523SKP\_SUP\_DONE (30,150). All...
- Highest **vocabulary density**: T19120722GS\_SUP\_DONE (127.1), T19110523SKP\_SUP\_DONE (116.6). Lowest density: T19120227NW\_SUP\_DONE (62.8), T19111205NW\_SUP\_DONE (63.0). All...
- Most **frequent words** in the corpus: said (20,303), prisoner (17,443), mr (10,032), went (8,499), saw (7,961). More...
- Words with **notable peaks in frequency** across the corpus: will ( ), day ( ), they ( ), wanted ( ), november ( ). More...
- Distinctive words** (compared to the rest of the corpus)
  - T19110717SKP\_SUP\_DONE: muller (74), brooks (66), mackenzie (54), coltman (52), stevens (53). More...
  - T19120319eg\_SUP\_DONE: ullman (61), gallafent (57), stevens

Words within each Document		
Document	Count	Difference
Type: prisoner		
27) T191 712		-7.64
3) T1912 642		-7.02
4) T1911 614		-6.05
40) T191 592		-6.80
19) T191 564		-6.89
9) T1912 554		-7.47
26) T191 547		-7.92
33) T191 545		-8.62
5) T1912 535		-6.46
21) T191 529		-8.69
7) T1912 504		-7.46
32) T191 500		-7.53

Keywords in Context		
Left	Keyword	Right
Document: T19101011CLR_SUP_DONE.xml		
Guardian" Newspaper Company.191010110006	Prisoner	had been twice convicted. He
that system.The Recorder sentenced	prisoner	to 20 months' hard labour, warn
Frederick Cook and Co. Each	prisoner	confessed to a previous convic
.d., with intent to defraud	Prisoner	bore an excellent character, an
,goods of N. Thierry, Limited	Prisoner	confessed to a previous convic
previous convictions were proved against	prisoner	(not for coining offences) dating
previous convictions were proved against	prisoner	(not for coining offences
had obtained employment.Sentences: Each	prisoner	, 15 months' hard labour; reco



## Words in the Entire Corpus

Word	Count	Trend
<input checked="" type="checkbox"/> election	18	
<input type="checkbox"/> selections	2	
<input checked="" type="checkbox"/> electioneering	2	

## Summary Corpus

- There are 41 documents in this corpus with a total of 2,535,564 words and 35,475 unique words.
- Longest documents** (by words): T19101011CLR\_SUP\_DONE (111,217), T19120227NW\_SUP\_DONE (103,304), Shortest documents: T19120722GS\_SUP\_DONE (22,965), T19110523SKP\_SUP\_DONE (30,150). All...
- Highest **vocabulary density** (by words): T19120722GS\_SUP\_DONE (127.1), T19110523SKP\_SUP\_DONE (116.6). Lowest density: T19120227NW\_SUP\_DONE (62.8), T19111205NW\_SUP\_DONE (63.0). All...
- Most frequent words** in the corpus: said (20,303), prisoner (17,443), mr (10,032), went (8,499), saw (7,961). More...
- Words with **notable peaks in frequency** across the corpus: will (~), day (~), they (~), wanted (~), november (~). More...
- Distinctive words** (compared to the rest of the corpus)
  - T19107175KP\_SUP\_DONE: mullet (74), brooks (66), mackenzie (54), coltman (52), stevens (53). More...
  - T19120319eg\_SUP\_DONE: ullman (61), gallafent (57), stevens (62), pearson's (52), hayes (53). More...
  - T19120109NW\_SUP\_DONE: minini (39), martin (57), rothwell (35), follett (30), clifford (31). More...
  - T19110905GS\_SUP\_DONE: august (156), fletcher (57), deceased (95), cabin (43), searle (37). More...
  - T19120611SKP\_SUP\_DONE: woods (100), stone (62), hole (53), wallet (44), titheradge (33). More...
 Next 5 of 36 remaining

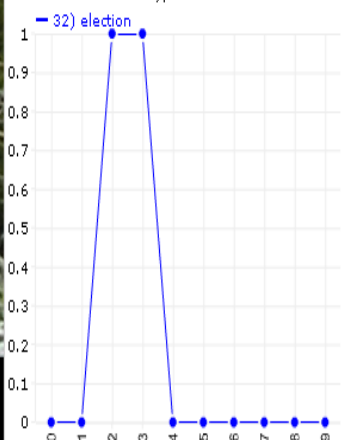
## Words within each Document

Document	Count	Difference	Relative	Trend
Type: election				
37) T191	7	-0.02	1.62	
40) T191	5	-0.05	0.74	
32) T191	2	-0.08	0.30	
35) T191	1	-0.09	0.15	
21) T191	1	-0.07	0.10	
8) T1911	1	-0.09	0.13	
3) T1912	1	-0.08	0.12	
41) T191	0	-0.10	0.00	
39) T191	0	-0.11	0.00	

Type: electioneering				
40) T191	1	-0.07	0.15	
29) T191	1	-0.08	0.16	
41) T191	0	-0.09	0.00	
39) T191	0	-0.10	0.00	
38) T191	0	-0.09	0.00	

## Word Trends

Distribution of Types across Document



## Keywords in Context Collocates

Collocate	Raw Frequency	Ratio
course	1 (8)	6.2
liberal	1 (4)	6.2
circumstances	1 (2)	6.2
despatched	1 (1)	6.2
candidate	1 (1)	6.2
barker	1 (1)	6.2
mr	2 (313)	12.0
george	2 (88)	12.4
lloyd	2 (25)	12.5

## *Contact Details*

- Bob Shoemaker ([r.shoemaker@shef.ac.uk](mailto:r.shoemaker@shef.ac.uk))
- Tim Hitchcock ([t.hitchcock@herts.ac.uk](mailto:t.hitchcock@herts.ac.uk))
- Michael Pidd (HRI) ([m.pidd@shef.ac.uk](mailto:m.pidd@shef.ac.uk))

## *Websites*

[www.oldbaileyonline.org](http://www.oldbaileyonline.org)

[www.londonlives.org](http://www.londonlives.org)

[www.connectedhistories.org](http://www.connectedhistories.org)

<http://criminalintent.org/>