

Web based computer-assisted text analysis

Christopher Ohge, Martin Steer

Riga Technical University, September 2019

Hathi Trust Research Centre (HTRC)



Hathi Trust Research Centre (HTRC)

- Promotes scholarly use of the HathiTrust Digital Library
- Custodian of Google books
 - OCR corpus – many character errors!
 - Biased – e.g C19th lacks genre metadata
- Provides
 - Data analysis infrastructure
 - “Non-consumptive” use
 - Extracted Features Dataset
 - Visualisation tools
 - Lots more...



HTRC Extracted Features Dataset

- Pre-extracted features
- Python/R libraries, API access
- Volume-level
 - Bibliographic metadata
 - title, pubDate, language,
 - genre, typeOfResource
 - names, imprint, rightsAttributes
 - pageCount, size

Data Stats

# of volumes	15,722,079
# of pages	5,787,519,444
# of tokens	2,449,739,213,773
# of IC volumes	9,914,509
# of IC pages	3,005,448,348
# of IC tokens	1,777,793,828,310
# of PD volumes	5,807,570
# pd pages	2,602,212,586
# pd tokens	1,197,838,539,662

HTRC Extracted Features Dataset

- Page level
 - Header, Body, Footer info
 - Page sequence numbers
 - tokenCount, lineCount, sentenceCount
 - Languages (inferred per page)
 - tokenPosCount (POS tags)
 - beginCharCounts, endCharCounts

```
{  
  "id": "uc1.b2617728",  
  "metadata": {  
    "schemaVersion": "1.3",  
    "dateCreated": "2016-06-25T03:45:05.180",  
    "volumeIdentifier": "uc1.b2617728",  
    "accessProfile": "google",  
    "rightsAttributes": "ic",  
    "hathitrustRecordNumber": "10281140",  
    "enumerationChronology": " ",  
    "sourceInstitution": "UC",  
    "sourceInstitutionRecordNumber": ".b165",  
    "oclc": [  
      "70545428"  
    ],  
    "isbn": [],  
  }  
}
```

HTRC Extracted Features Dataset

- Page level
 - Header, Body, Footer info
 - Page sequence numbers
 - tokenCount, lineCount, sentenceCount
 - Languages (inferred per page)
 - tokenPosCount (POS tags)
 - beginCharCounts, endCharCounts

```
{  
  "seq": "00000033",  
  "tokenCount": 273,  
  "lineCount": 36,  
  "emptyLineCount": 0,  
  "sentenceCount": 11,  
  "header": {  
    "tokenCount": 0,  
    "lineCount": 0,  
    "emptyLineCount": 0,  
    "sentenceCount": 0,  
    "capAlphaSeq": 0,  
    "beginCharCounts": {},  
    "endCharCount": {},  
    "tokenPosCount": {}  
  },  
  "body": {  
    "tokenCount": 273,  
    "lineCount": 36,  
    "emptyLineCount": 0,  
    "sentenceCount": 11,  
    "capAlphaSeq": 0,  
    "beginCharCounts": {},  
    "endCharCount": {},  
    "tokenPosCount": {}  
  }  
}
```

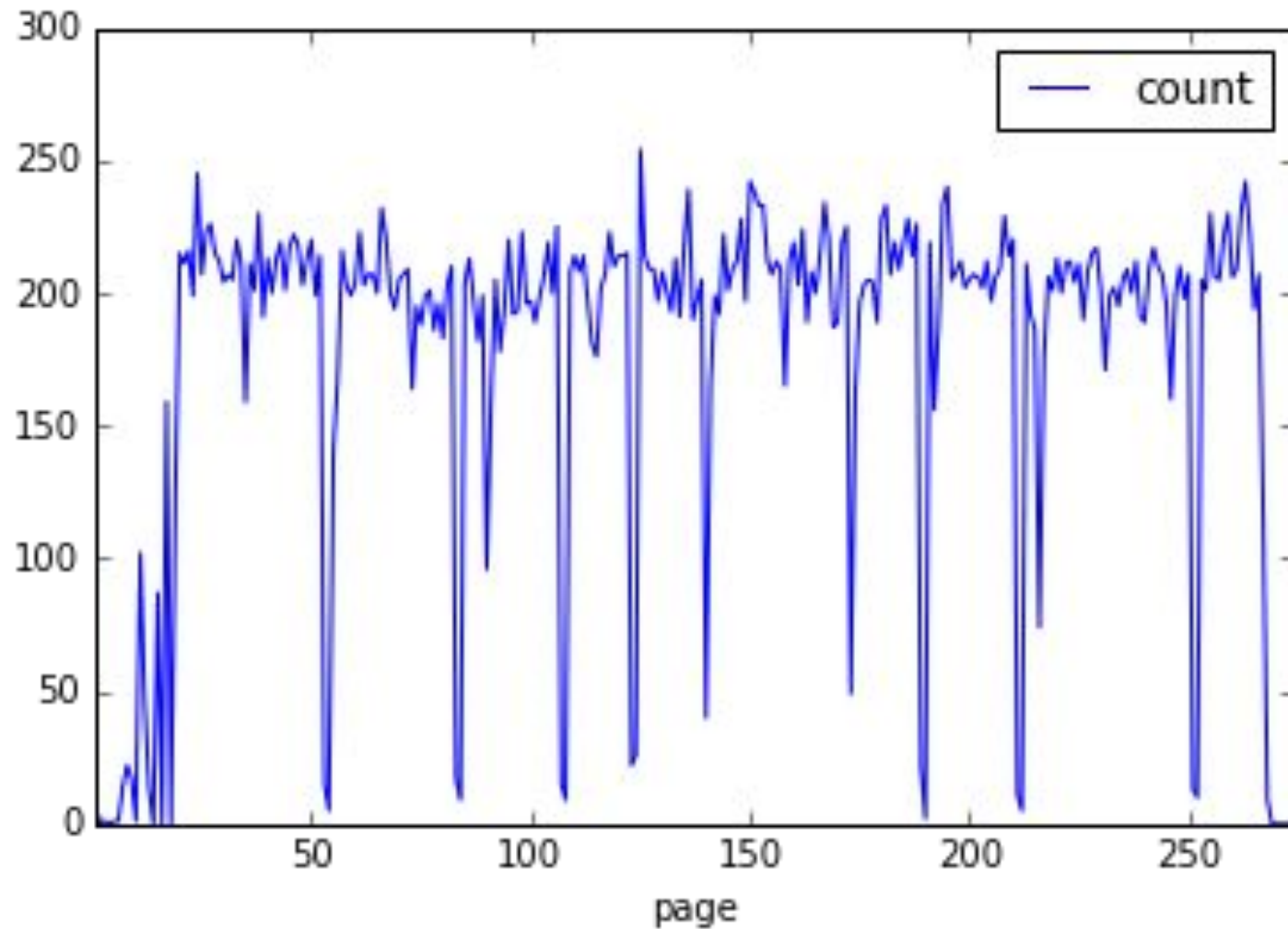
HTRC Extracted Features - Dataframe

page	section	token	pos	count
27	body	those	DT	1
		within	IN	1
28	body	a	DT	3
		be	VB	1
		deserted	VBN	1
		faintly	RB	1
		important	JJ	1

HTRC Extracted Features - Dataframe

section	token	count
body	,	3258
	"	1670
	the	1565
	.	1532
	and	1252

HTRC Extracted Features - Plotting

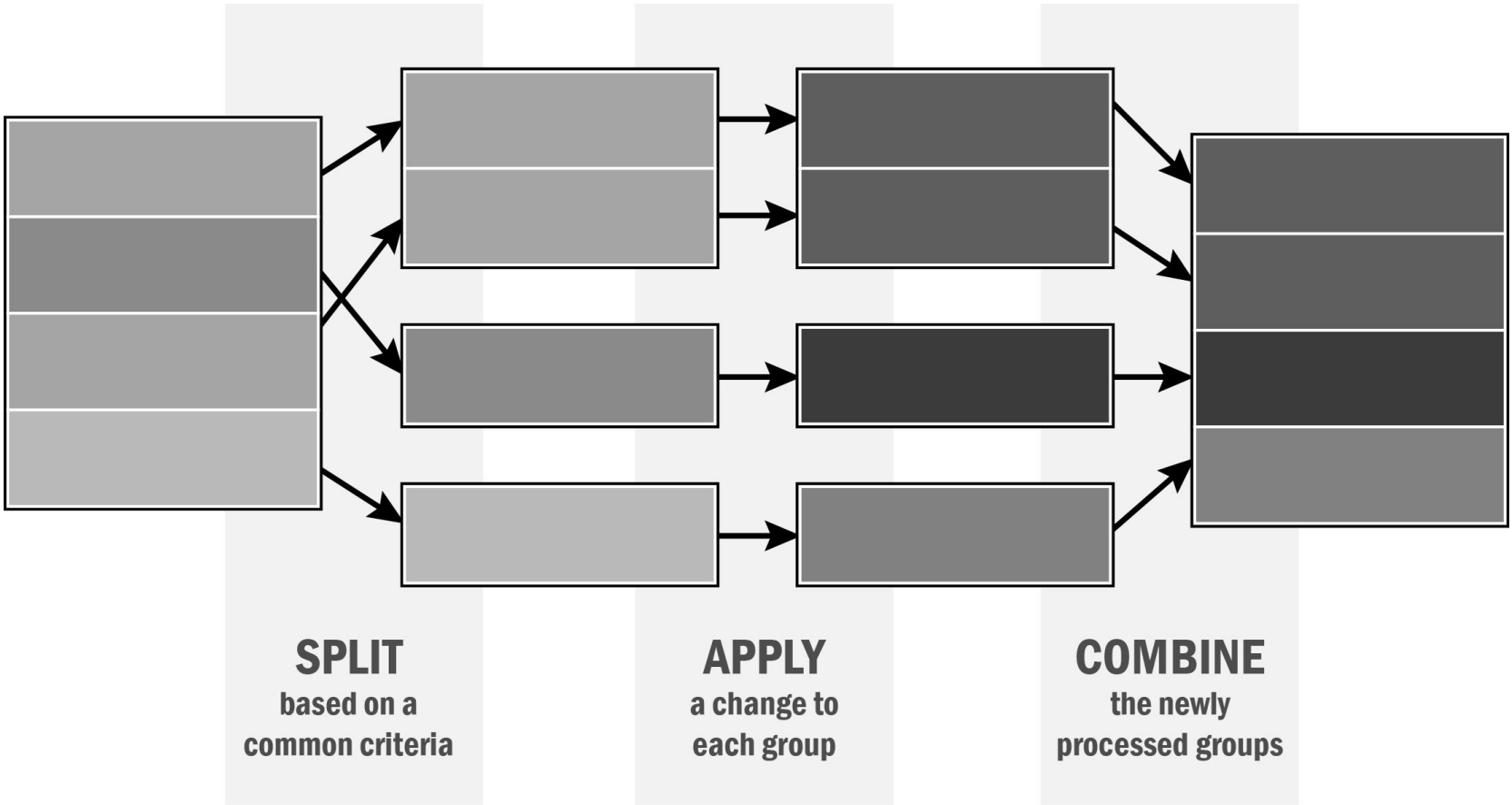


HTRC Extracted Features - Grouping

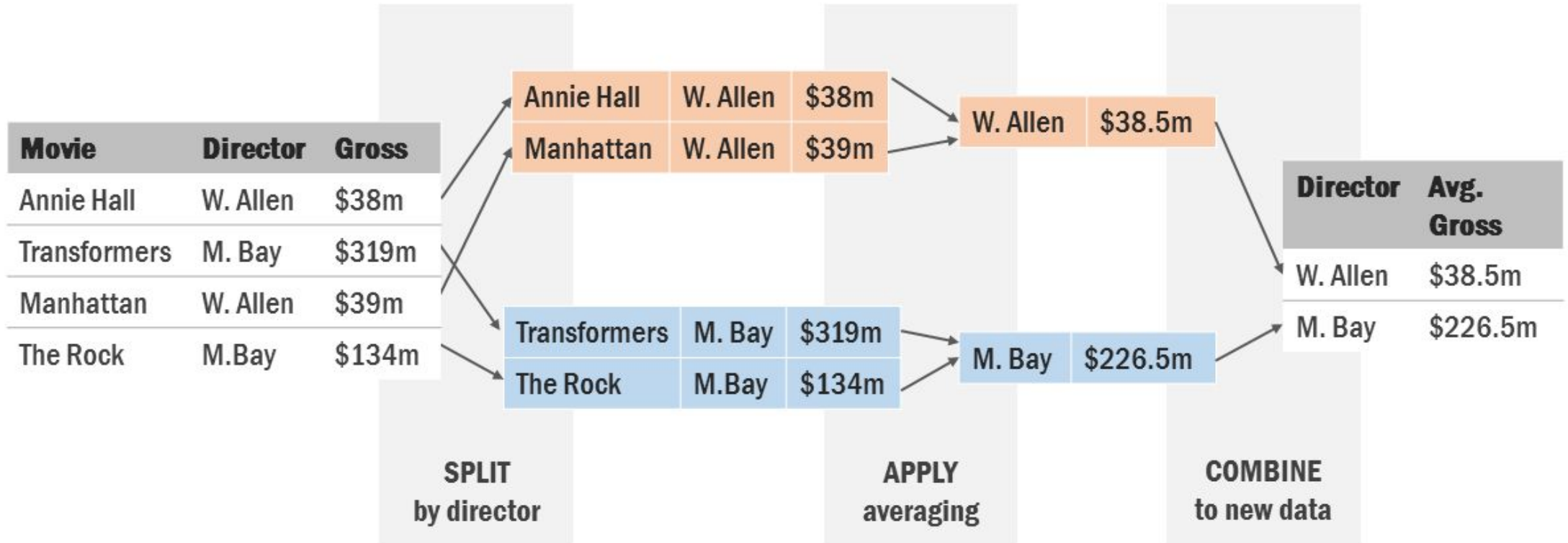
page	section	token	pos	count
19	header	CHAPTER	NNP	1
35	header	CHAPTER	NNP	1
56	header	CHAPTER	NNP	1
73	header	CHAPTER	NNP	1
91	header	CHAPTER	NNP	1
115	header	CHAPTER	NNP	1
141	header	CHAPTER	NNP	1
158	header	CHAPTER	NNP	1
174	header	CHAPTER	NNP	1
193	header	CHAPTER	NNP	1

HTRC Extracted Features - Plotting

Split-apply-combine

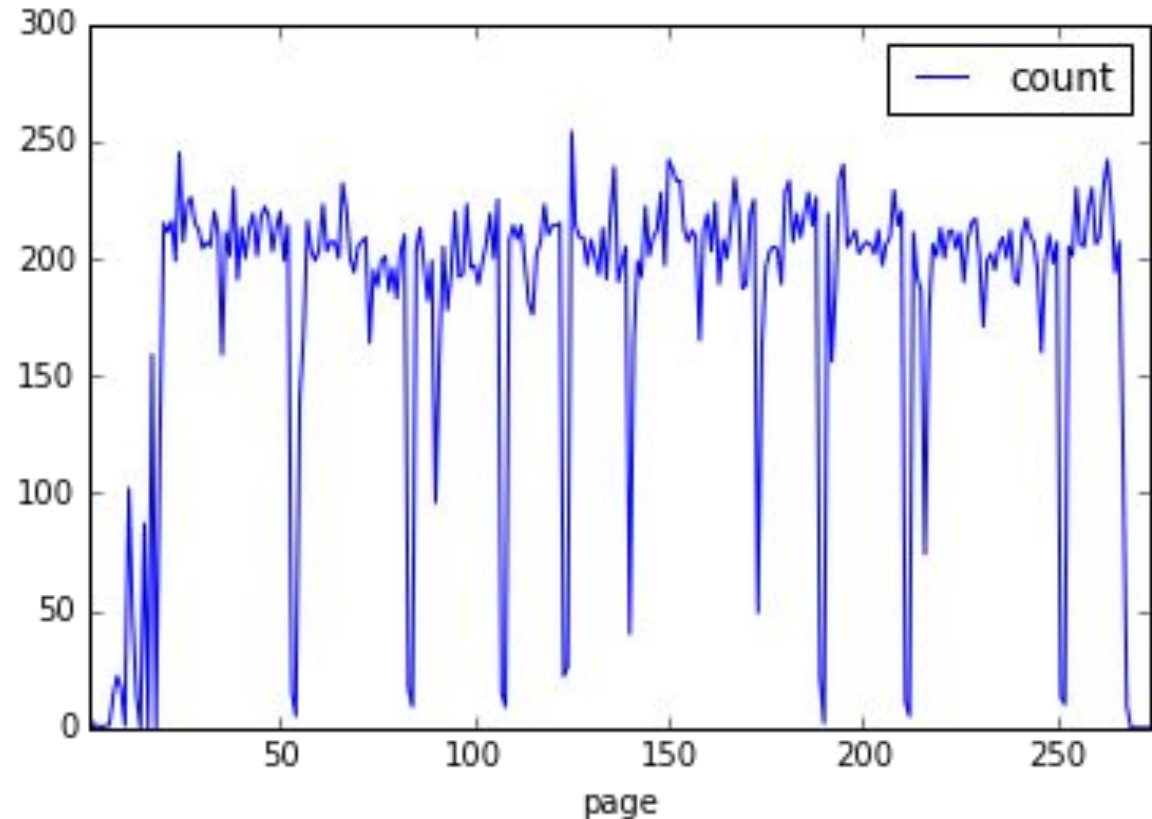


Split-apply-combine



Text analysis use cases

- Distant reading
- Word similarity
- Topic models
- Plot analysis
- Emotion analysis
- Visual structure



<https://programminghistorian.org/en/lessons/text-mining-with-extracted-features>

Text analysis use cases

- Distant reading
- Word similarity
- Topic models
- Plot analysis
- Emotion analysis
- Visual structure

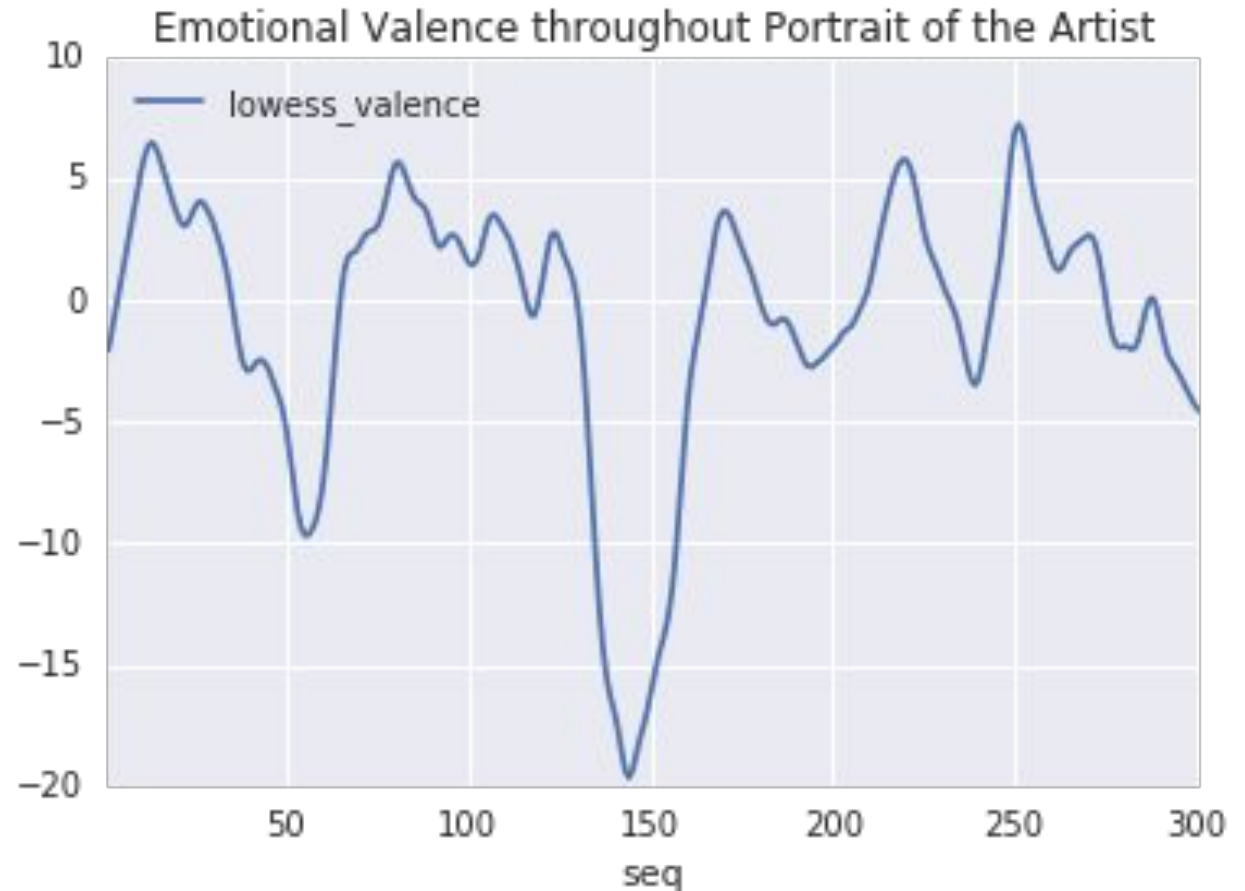
Search for:

1923	washington	columbia	ohio	york	chicago	montgomery	west	boston	michigan
1922	washington	dawes	pennsylvania	columbia	york	dent	president	city	washington.
1921	washington	columbia	washington.	national	pennsylvania	president	dent	boston	york
1920	washington	boston	columbia	illinois	york	washington.	ohio	city	pennsylvania
1919	washington	illinois	address	york	vice	national	boston	washington.	tribune
1918	washington	york	illinois	boston	baltimore	american	national	president	pennsylvania
1917	washington	boston	washington.	columbia	union	address	charleston	baltimore	vice
1916	washington	washington.	boston	virginia	columbia	charleston	massachusetts	wilson	president
1915	washington	washington.	charleston	virginia	massachusetts	louis	columbia	baltimore	portsmouth
1914	washington	columbia	denver	philadelphia	maryland	louis	baltimore	virginia	ohio
1913	washington	philadelphia	baltimore	charleston	american	address	pennsylvania	columbia	massachusetts
1912	washington	va.	columbia	charleston	louis	columbian	baltimore	boston	american

<https://wiki.htrc.illinois.edu/display/COM/Extracted+Features+in+the+Wild>

Text analysis use cases

- Distant reading
- Word similarity
- Topic models
- Plot analysis
- Emotion analysis
- Visual structure

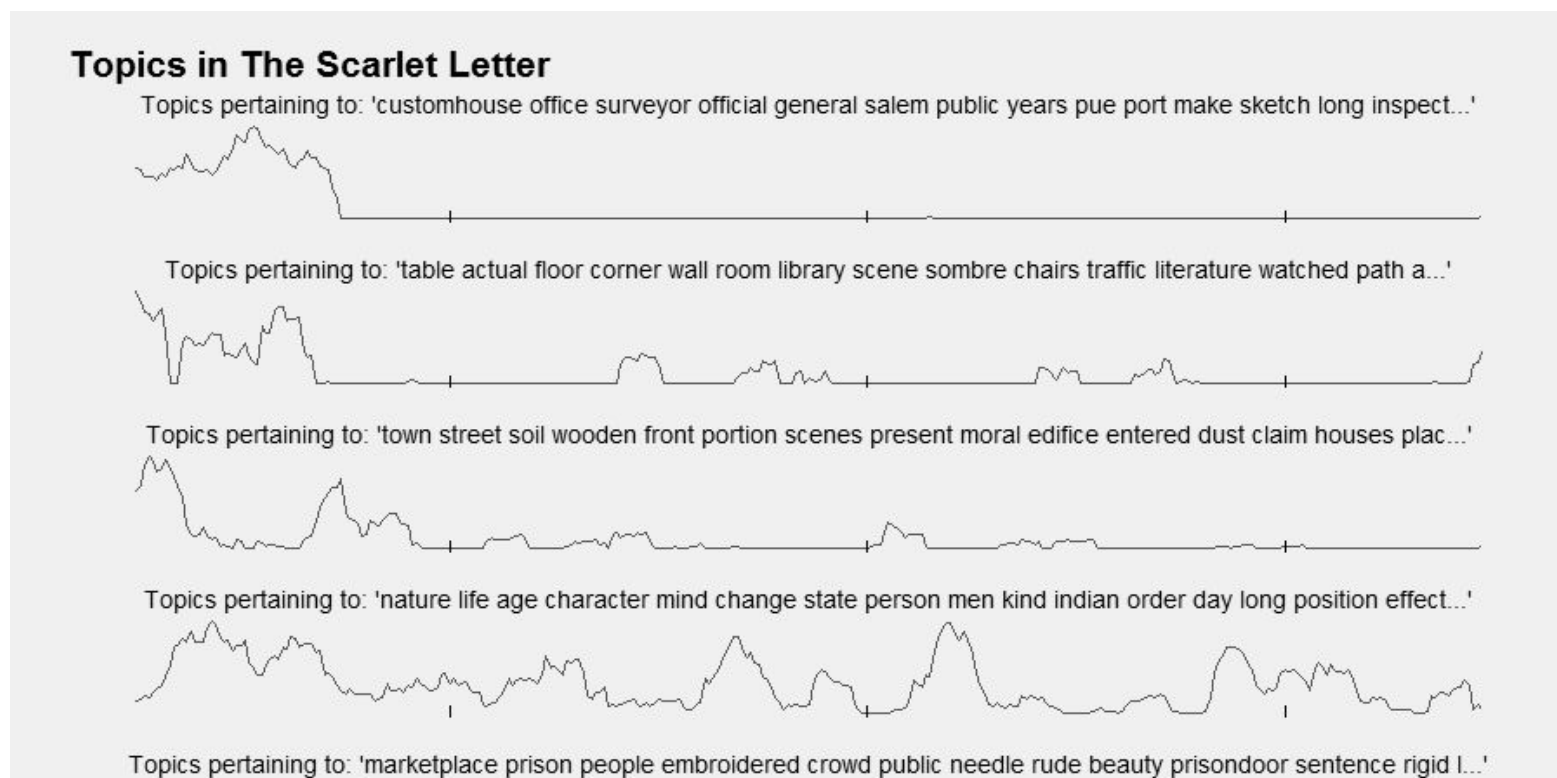


<https://programminghistorian.org/en/lessons/text-mining-with-extracted-features>

Text analysis use cases

- Distant reading
- Word similarity
- Topic models
- Plot analysis
- Emotion analysis
- Visual structure

Within-book topic modelling, Peter Organisciak

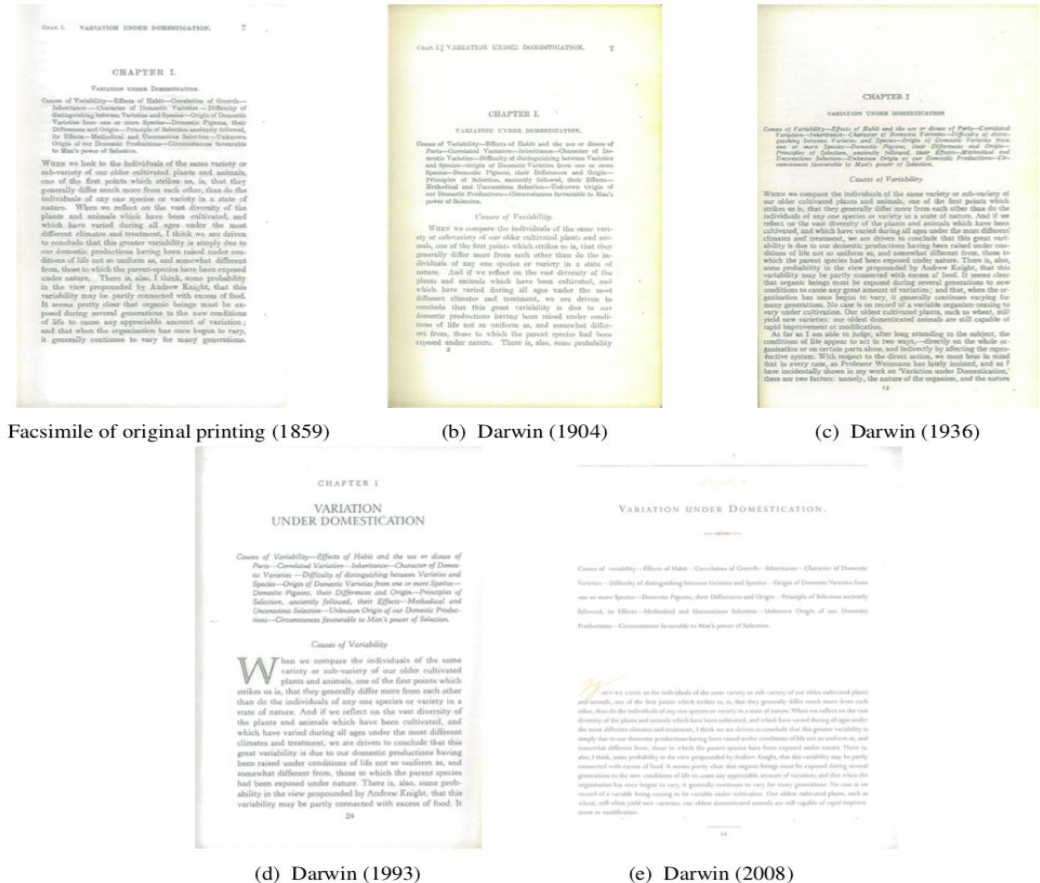


Interpretation of mediums

5 copies of Origin of Species in HTC corpora

Two trends:

- 1929 and after are on average about one quarter of an inch taller and one eighth of an inch wider than books published before
- the font size (as determined by the median line height) also appears to be increasing (by a quarter of a point)



(a) Facsimile of original printing (1859) (b) Darwin (1904) (c) Darwin (1936)

Figure 1: Five versions of Darwin's *Origin of Species*, illustrating different design choice in page layout.

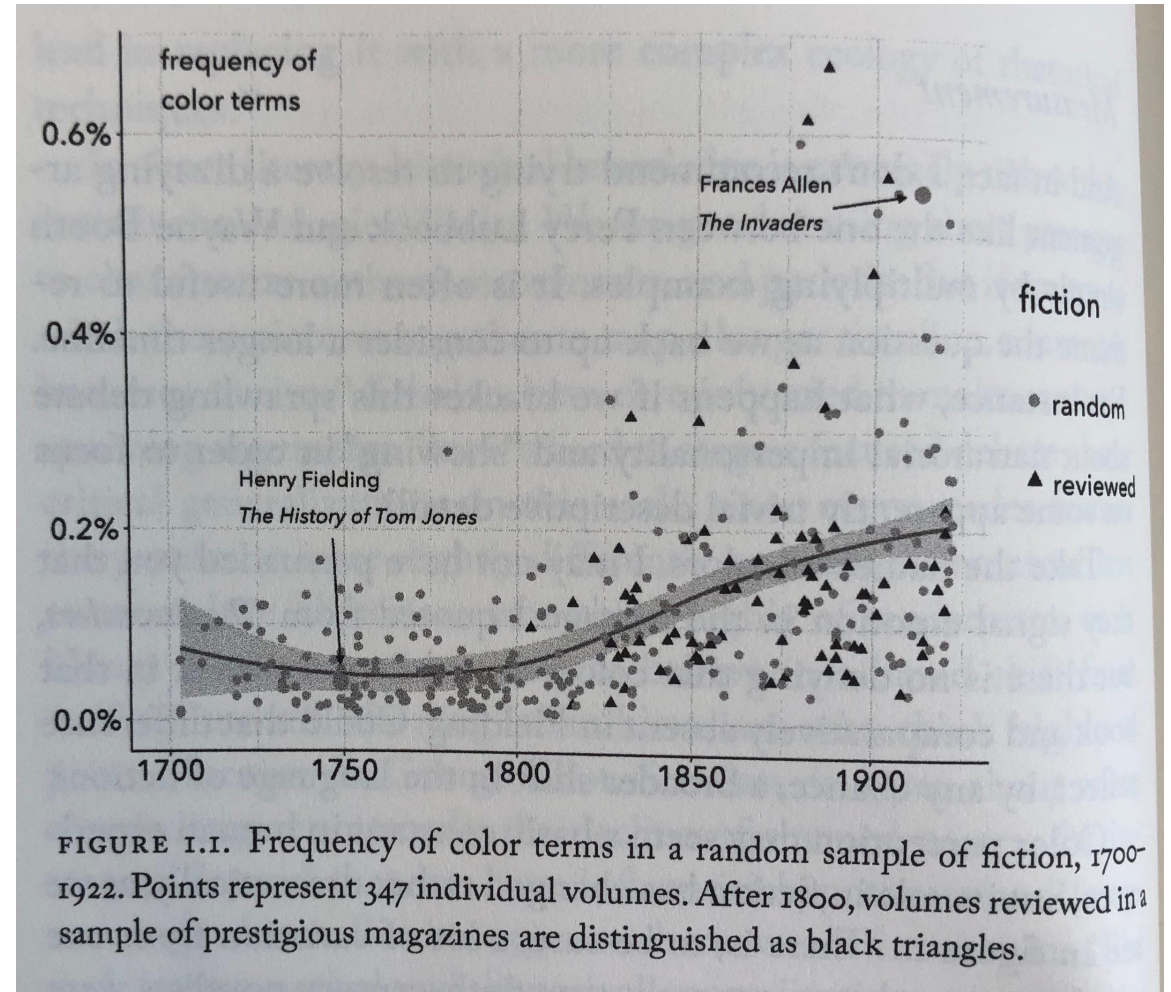
Bamman, David, and Björn Hartmann. 'Modeling the History of Book Design HTRC Whitepaper: Summary of Activities', n.d., 8.

Insight driven investigation

Distant Horizons, Ted Underwood (2019)

Frequency of colour terms in random fiction increases post 1800 – why?

- Does it represent the decline in 3rd person narration? A shift to description?
- We need to look for correlations

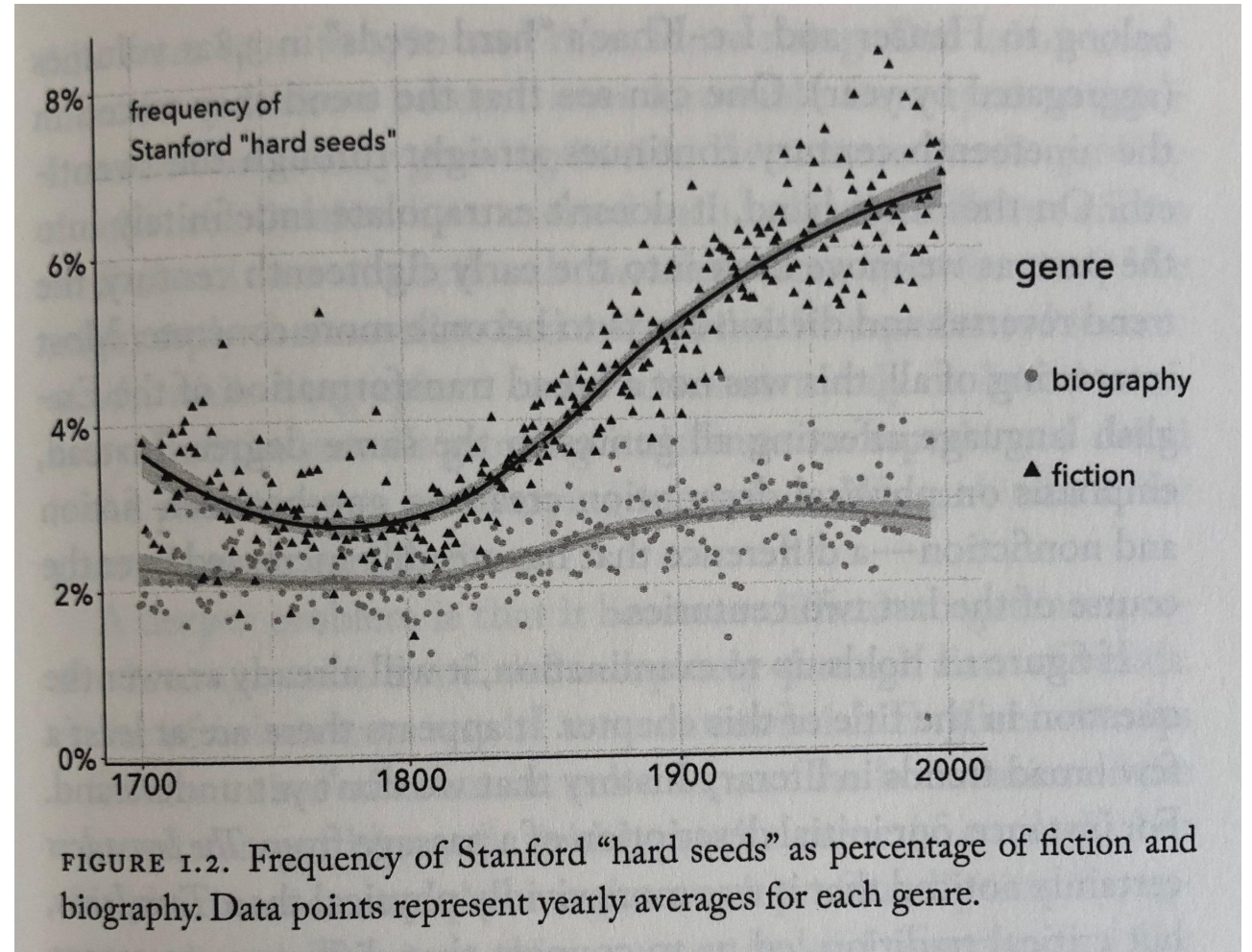


Frames of context

Standord 2012 Heuser & Le-Khac:

- Rising frequency of concrete adjectives in 19th century novels (colours, names, body parts, action verbs – used for physical descriptions)

Concluded fictional discourse was transitioning "from telling to showing"



Practical – Finding Trends

bookworm: [HathiTrust](#)

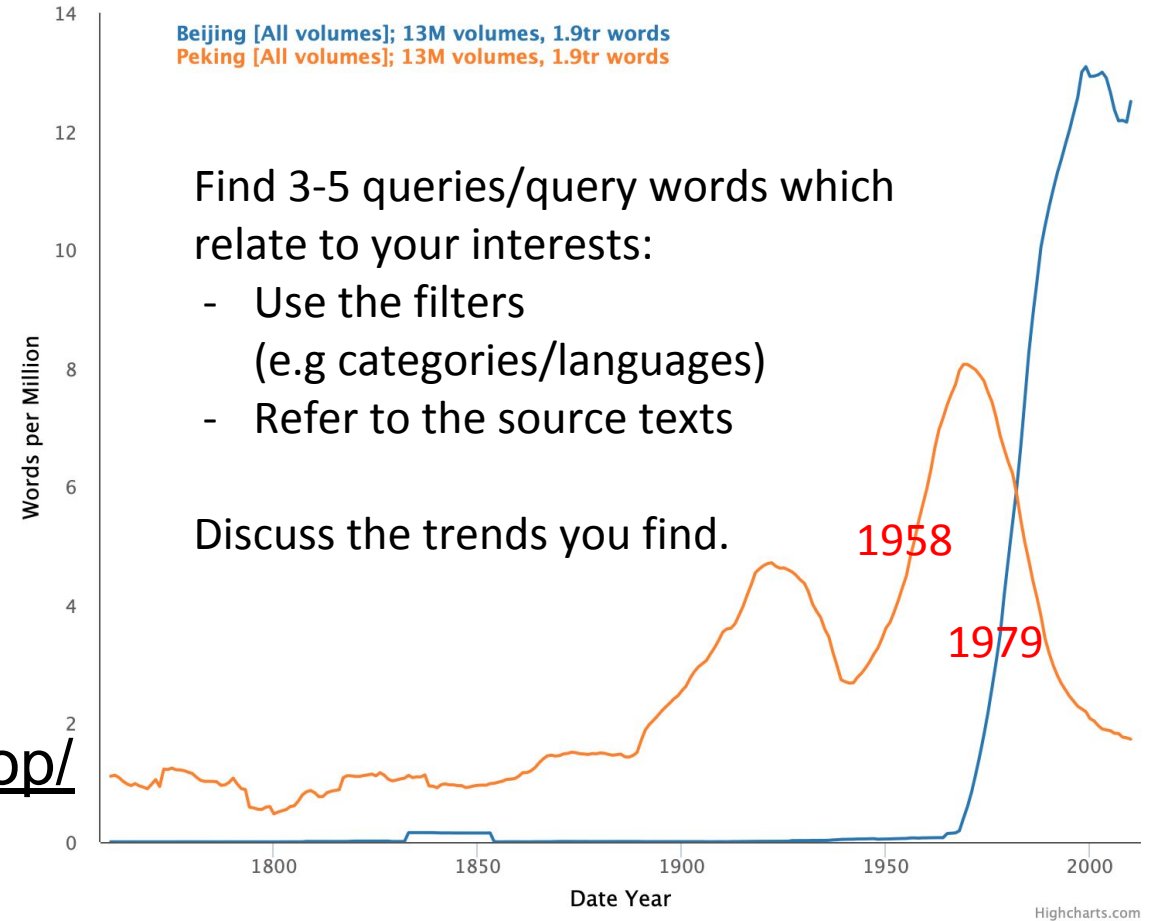
Search for trends in millions of volumes at <http://hathitrust.org>

Beijing in All volumes

Peking in All volumes

<https://bookworm.htrc.illinois.edu/develop/>

<https://books.google.com/ngrams>



Find 3-5 queries/query words which relate to your interests:

- Use the filters (e.g categories/languages)
- Refer to the source texts

Discuss the trends you find.

Examples

Literary trends

- Red
- White

(Narrow class filter:
fiction/biography)

Meaning change

- Hysteria
- Gay
- legacy/heritage

Parts of speech

Book_VERB
Book_NOUN

Geopolitical

- Persia
- Iran
- ایران

- Beijing
- China

Practical – Interpreting Topics

HathiTrust Fiction (1920-1922) topic-browser

<https://jgoodwin.net/htb/#/model/grid>



Pick 3 topics in each group:

- Use the GUI to investigate
- Why does the GUI help?
- Refer to the source texts

Discuss your findings.

French literature cluster?

Label = French Women or French Gender?

Contains = Fiction with a lot of gender descriptions?

A lot of polite conversation with women characters?

Computer-assisted tools and Hermeneutics

Descartes, *Discourse on Method*: “[T]hings made up of different elements and produced by the hands of several master craftsmen are often less perfect than those on which only one person has worked.”

Individual versus communal-participatory inquiry; dialogic relationship between human and human and machine, criticism and code; augmentation of criticism.

Teamwork changes the single-interpreter model that Descartes describes: “analytical tools *are* instantiations of interpretive methods that can be woven closely into other hermeneutical things” (Rockwell and Sinclair).