
Text Analysis with Voyant Tools & AntConc

Mia Partlow and Erica Hayes
NCSU Libraries

Learning Goals

- Introduce you to Text Analysis & best practices for building a corpus
- Learn the capabilities of Voyant Tools and AntConc for text reading, analysis and visualization
- Load files of various types and URLs into Voyant Tools & AntConc
- Navigate the interfaces of both tools in multiple ways
- Export your views and download images of views

What is Text Analysis?

“Text Analysis is the search for and discovery of patterns and trends in a corpus of texts. The analysis of those patterns and trends can help researchers uncover previously unseen characteristics of a specific corpus, deconstruct a text, and reveal new ideas and theories about a particular genre, topic or author.” -Meredith Dabek (2014)

Two guiding questions for text analysis -Sinclair & Rockwell (2015):

- For texts with which I am already familiar, how can computers help me identify and study interesting things I had not noticed before, or things I had noticed but did not have reasonable means to pursue?
- How can computers help me identify and understand texts with which I am not familiar with or which I cannot reasonably read?

Concepts & Challenges

- **Corpus:** a collection of texts.
 - Consider the size of your corpus and your inclusion and exclusion criteria
- **Data cleaning:** removing information not relevant to your analysis.
 - Examples: chapter numbers and titles, copyright information, headers and footers, page numbers.
 - Languages and tools such as R, Python, and Regular Expressions are necessary with a large corpus.
- **Stopwords:** words you filter out of your analysis
 - Standard lists exist (e.g. [Buckley-Salton Stoplist](#))
 - Voyant has a stoplist included
 - Customization of the stoplist may be necessary

Choosing a Corpus

Text Data Sources:

- [Project Gutenberg ebooks](#)
- [HathiTrust Digital Library](#)
- [Datasets and Scholarly APIs at NC State](#)
- Surveys and interview data

“Quantitative analysis tends to require context before it becomes meaningful. It doesn’t mean much to say that the word “motion” is common in Wordsworth, for instance, until we know whether “motion” is more common in his works than in other nineteenth-century poets. So yes, text-mining can provide clues that lead to real insights about a single author or text. But it’s likely that you’ll need a collection of several hundred volumes, for comparison, before those clues become legible”

- Ted Underwood (2012)

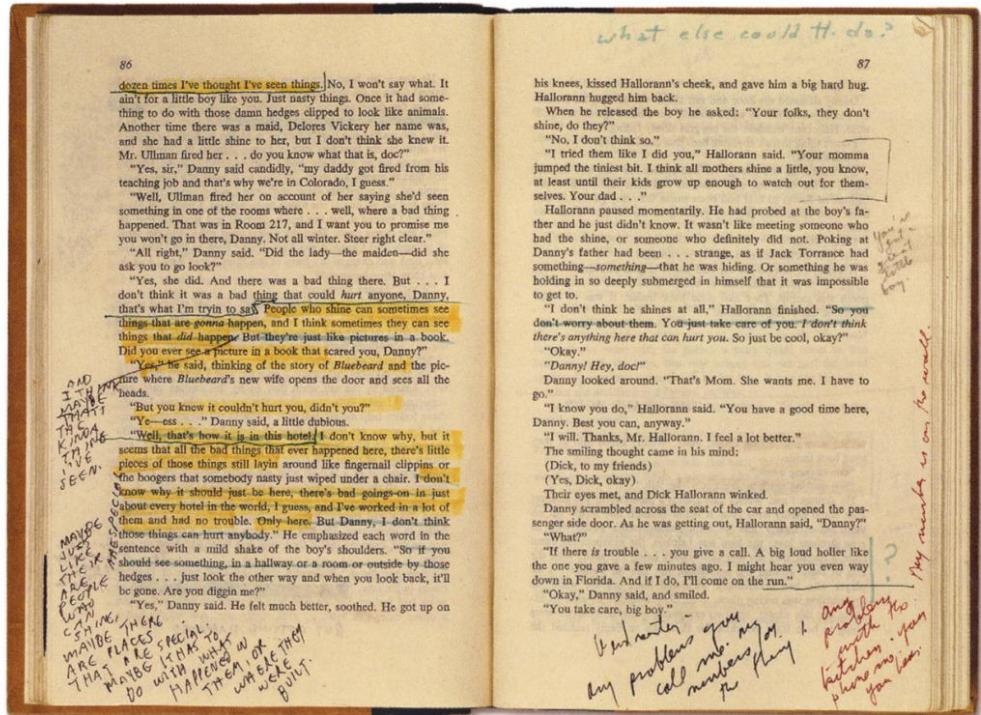
Copyright

- Mining databases or corpora to conduct scholarly analysis, without subsequently republishing the contents of those databases or corpora, can be fair use under [Author's Guild v. HathiTrust, 755 F.3d 87 \(2d Cir. 2014\)](#).
- A database or website you are using to create your corpus might regulate text mining, or might allow text mining but prohibit sharing excerpts from the materials that they made available for mining. Therefore, it's important to carefully read any terms of use or licenses as you prepare your project for publication.

Close Reading or Microanalysis

Close Reading/Microanalysis

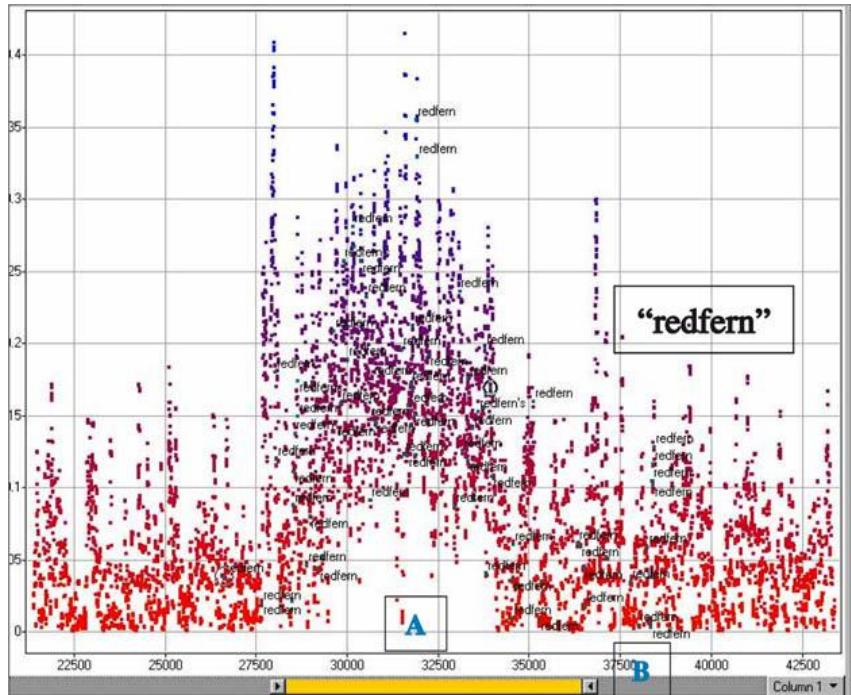
“Essentially, close reading means reading to uncover layers of meaning that lead to deep comprehension.” - Nancy Boyles (2013)



Distant Reading or Macroanalysis

Distant Reading/Macroanalysis

“Distant reading aims to generate an abstract view by shifting from observing textual content to visualizing global features of a single or of multiple text(s)” -Stefan Jänicke (2016)



Benefits of both forms of reading and analysis

“It is the exact interplay between the macro and micro scale that promises a new, enhanced, and perhaps even better understanding of the literary record. The two approaches work in tandem and inform each other. Human interpretation of the ‘data,’ whether it be mined at the macro or micro level, remains essential. While the methods of enquiry, of evidence gathering, are different, they are not antithetical, and they share the same ultimate goal of informing our understanding of the literary record, be it writ large or small”

-Matthew Jockers “On Distant Reading and Macroanalysis” (2011)

Text Analysis activity goals for today's workshop

- Explore word frequencies, occurrences of keywords, concordances and collocations in 8 fairy tales from Project Gutenberg
- The goal of this activity is to introduce you to the Voyant Tools & AntConc interfaces to enable quantitative analysis of the texts.

Voyant Tools

Voyant: The Good & The Bad

Positives:

Allows for quick reading and analysis of texts

Easier to identify themes and topics to explore further

Comparative options to analyze texts in tandem

~21 tools for visualizing the texts

HTML5 based tools (unlike Voyant 1.0)

Negatives:

Web application frequently breaks down/fails to load

Slow loading time with large documents (web)

Interface limits number of texts you can analyze

Doesn't allow you to see under the hood of the tools calculations

Voyant Server

Voyant Server is a version of Voyant that you can download and run locally on your computer. It's a Java-based application.

- Avoid problems of slow loading time and instability of web version
- Keep data private

[Installation](#) instructions

[Latest release download](#) (must have [Java downloaded](#))



voyant-tools.org

Add Texts ?

Type in one or more URLs on separate lines or paste in a full text.

Open Upload Reveal

Voyant Tools is a web-based reading and analysis environment for digital texts.

Accepted file types

File types: plain text, HTML, XML, PDF, RTF, and MS Word

Zipped directories OK - (.zip, .tar, .tgz, etc.) containing documents in those formats.

Get the data

go.ncsu.edu/voyantfall2018

Download the “**fairy-tales-corpus.zip**” file

**Can't access the link? Make sure you are logged into your NCSU gmail!

Getting started

Choose either 1 or 2:

1. Web-based version:

Open a web browser, visit <http://voyant-tools.org>

2. Voyant Server (highly recommended!):

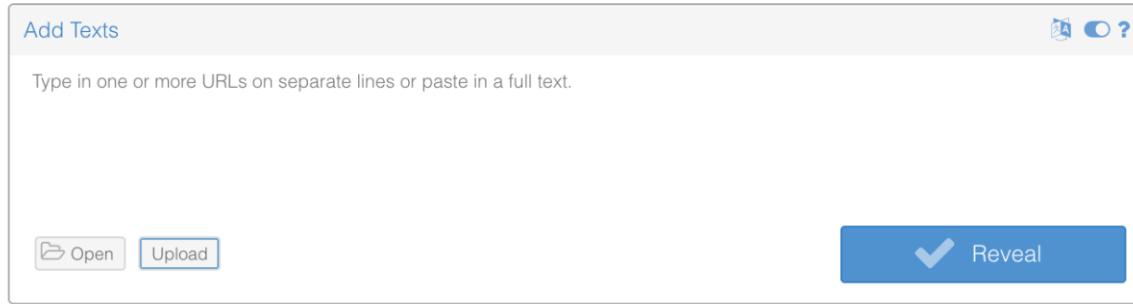
- [download the VoyantServer 2.0 zip archive](#)
- double-click on the zip archive to expand its contents
- double-click on VoyantServer.jar
 - on Mac, because of security restrictions on applications that aren't signed and approved by Apple, you may need to Ctrl-click on the VoyantServer.jar file, select open from the menu, and then click open (not the default button) in the next dialog box

Load data into Voyant

The data you load into Voyant is called a “corpus.” Three ways to load a corpus:

1. Paste text into the large box
2. Paste URLs into the large box, one URL per line
3. Click the upload button and select files to upload

Uploading corpus



Voyant Tools is a web-based reading and analysis environment for digital texts.

Do Now:

Click Upload button and upload **fairy-tales-corpus.zip**

Tab through different tools



Scale Terms:

 Summary  Documents  Phrases

This corpus has 8 documents with 27,732 total words and 2,999 unique word forms. Created now.

Document Length:

- Longest: **THE LITTLE MERMAID** (9223); **BEAUTY AND THE BEAST** (5733); **THUMBELINA** (4381); **LITTLE SNOW WHITE** (2290); **CINNAMON** (1990)
 - Shortest: **THE FROG-PRINCE** (1199); **RAPUNZEL** (1395); **BRIAR ROSE** (1514); **CINDERELLA** (1997); **LITTLE SNOW WHITE** (2290)

Vocabulary Density:

- Highest: RAPUNZEL (0.303); BRIAR_ROSE (0.301); THE_FROG-PRINCE (0.284); LITTLE_SNOW_WHITE (0.244); CINDERELLA (0.244)
 - Lowest: THE_LITTLE_MERMAID (0.171); THUMBELINA (0.201); BEAUTY AND THE BEAST (0.202); CINDERELLA (0.223); LITTLE

Average Words Per Sentence:

- Highest: **BRIAR ROSE** (36.0); **THE FROG-PRINCE** (29.2); **THE LITTLE MERMAID** (27.8); **CINDERELLA** (25.9); **RAPUNZEL** (25.8)
 - Lowest: **THUMBELINA** (23.2); **LITTLE SNOW WHITE** (24.4); **BEAUTY AND THE BEAST** (25.5); **RAPUNZEL** (25.8); **CINDERELLA** (25.9)

Most frequent words in the corpus: said (201); little (152); came (89); beautiful (81); prince (74)

Reader TermsBerry

THUMBELINA

There was once a woman who wished very much to have a little child, but she could not obtain her wish. At last she went to a fairy, and said, "I should so very much like to have a little child; can you tell me where I can find one?"

"Oh, that can be easily managed," said the fairy. "Here is a barleycorn of a different kind to those which grow in the farmer's fields, and which the chickens eat; put it into a flower-pot, and see what will happen."

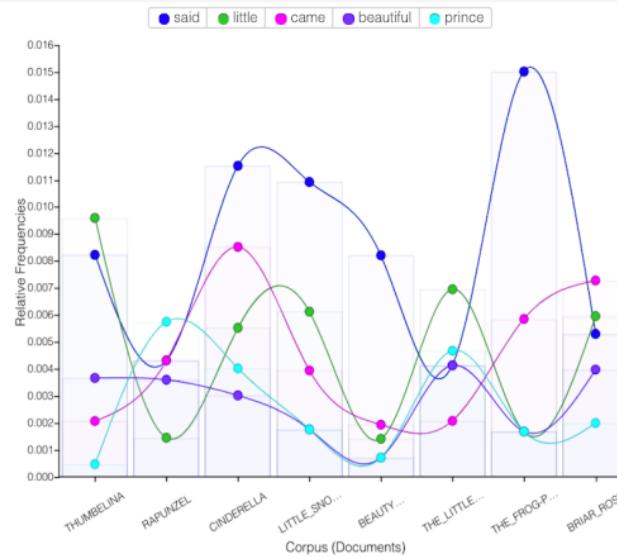
"Thank you," said the woman, and she gave the fairy twelve shillings, which was the price of the barleycorn. Then she went home and planted it, and immediately there grew up a large handsome flower, something like a tulip in appearance, but with its leaves tightly closed as if it were still a bud. "It is a beautiful flower," said the woman, and she kissed the red and golden-colored leaves, and while she did so the flower opened, and she could see that it was a real tulip. Within the flower, upon the green velvet stamens, sat a very delicate and graceful little maiden. She was scarcely half as long as a thumb, and they gave her the name of "Thumbelina," or Tiny, because she was so small. A walnut-shell, elegantly polished, served her for a cradle; her bed was formed of blue violet-leaves, with a rose-leaf for a counterpane. Here she slept at night, but during the day she amused herself on a table, where the woman had placed a number of plates. These plates were wreaths of flowers with their stems in the middle, and the tulip-leaf, which served Tiny as a chair, she would turn herself from side to side, with two oars made of white poppy-stems, singing prettily. Tiny could, also, sing so softly and sweetly that nothing like her singing had ever before been heard.

Select & Change Tools

Export

Help &
Info

Set Stopwords



Document		Left	Term	Right
□	1) THU...	went to a fairy, and	said	, "I should so very much
□	1) THU...	that can be easily managed,"	said	the fairy. "Here is a
□	1) THU...	what will happen." "Thank you,"	said	the woman, and she gave
	1) THU...	It is a beautiful flower,"	said	the woman, and she kissed
	1) THU...	would make for my son,"	said	the toad, and she took
	1) THU...	loud, or she will wake,"	said	the toad, "and then she
□	1) THU...	her in the water, and	said	, "Here is my son, he
□	1) THU...	toad, and heard what she	said	, so they lifted their heads
□	1) THU...	turned up their feelers, and	said	, "She has only two legs
□	1) THU...	looks." "She has no feelers,"	said	another. "Her waist is quite
□	1) THU...	being." "Oh! she is ugly."	said	all the lady cockchafers, although

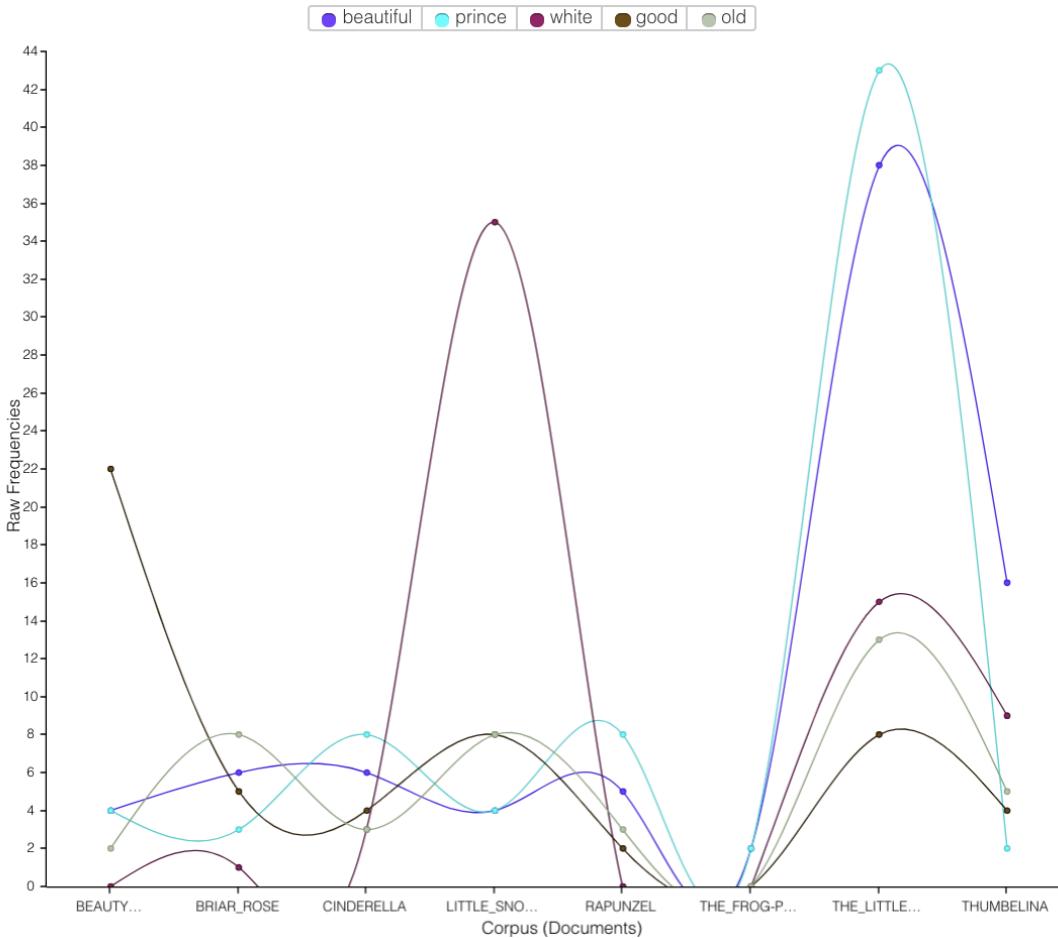
Start hands-on exploration of the corpus

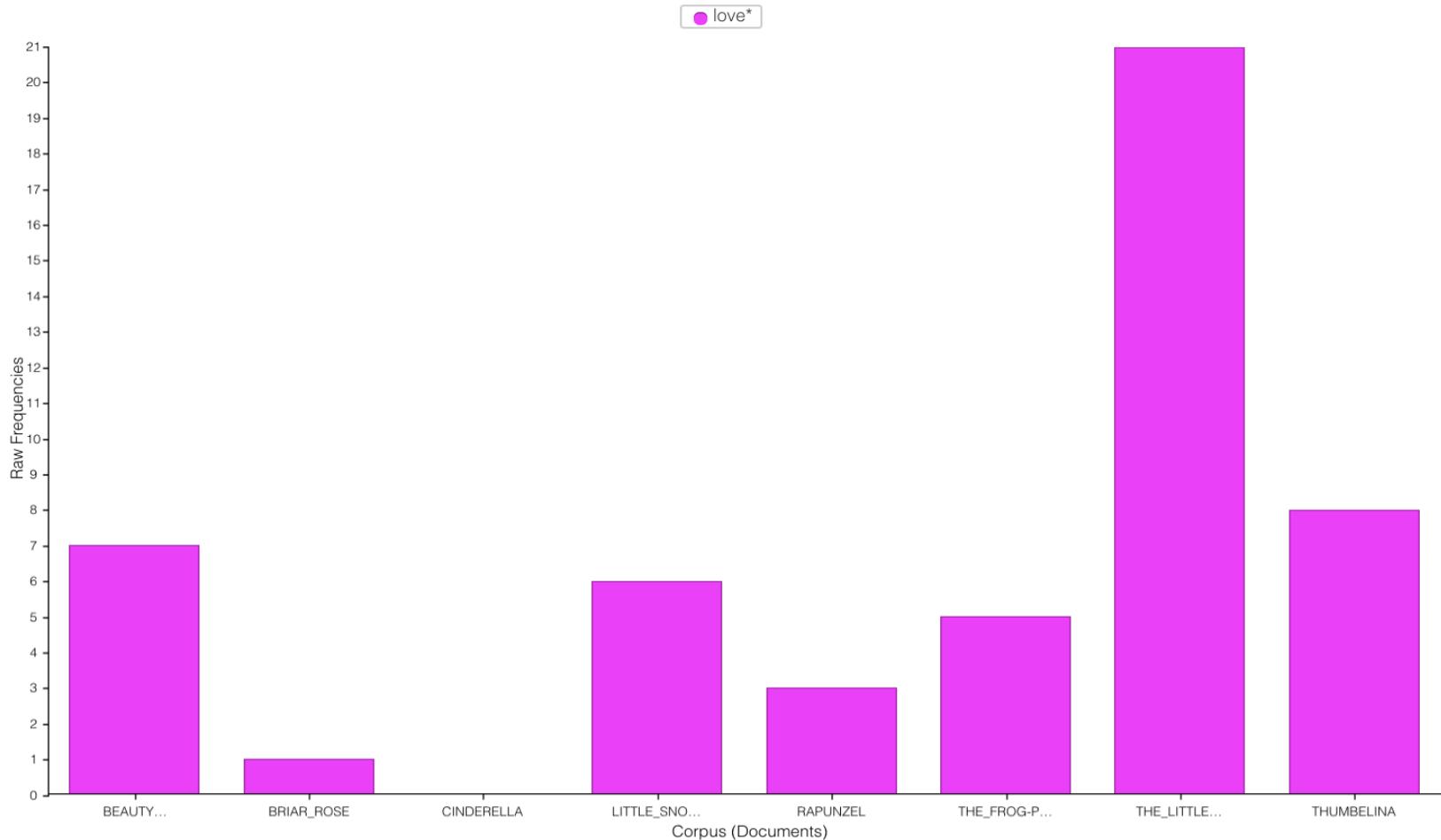
See Activity Guide handout

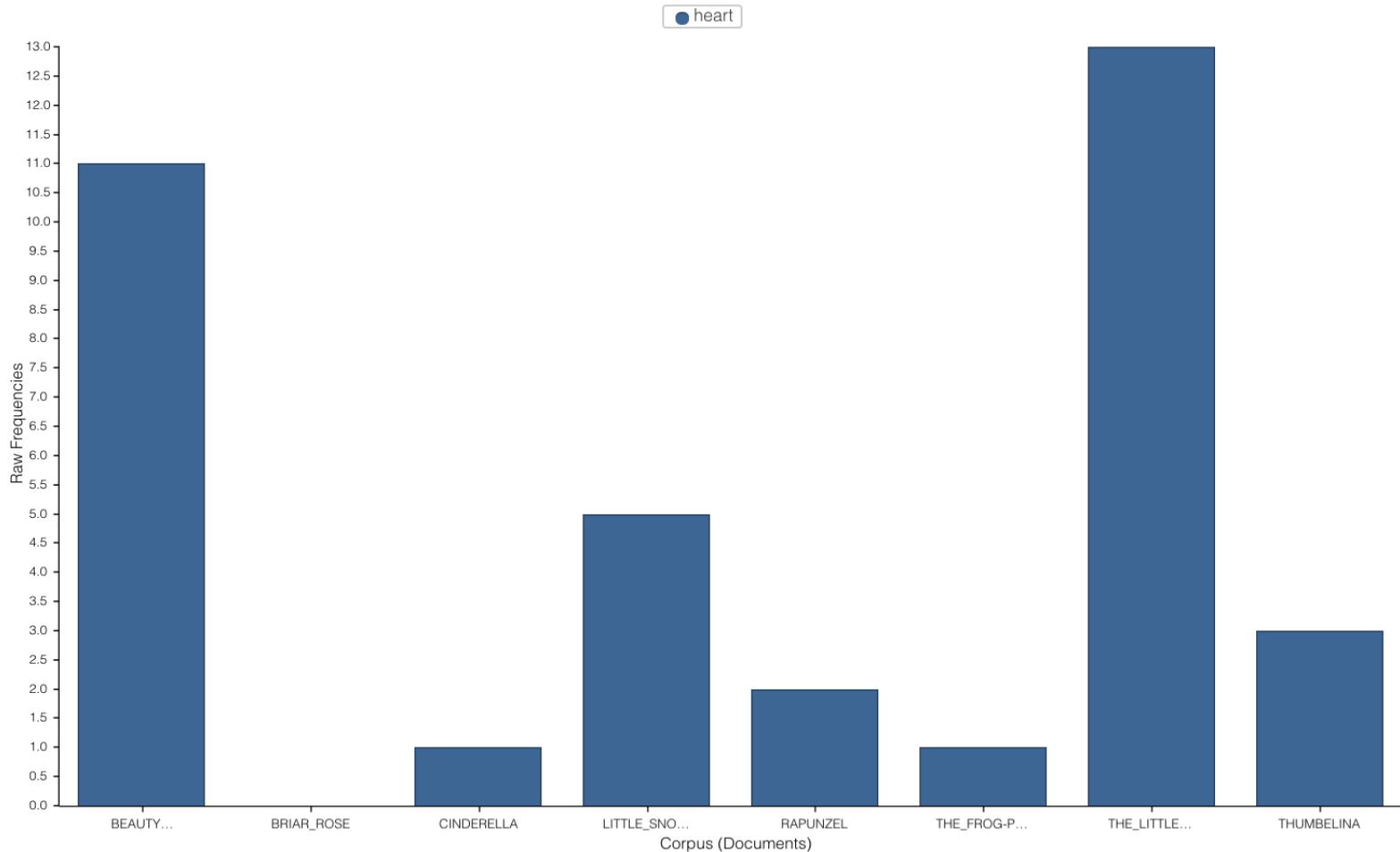
for self-directed activity

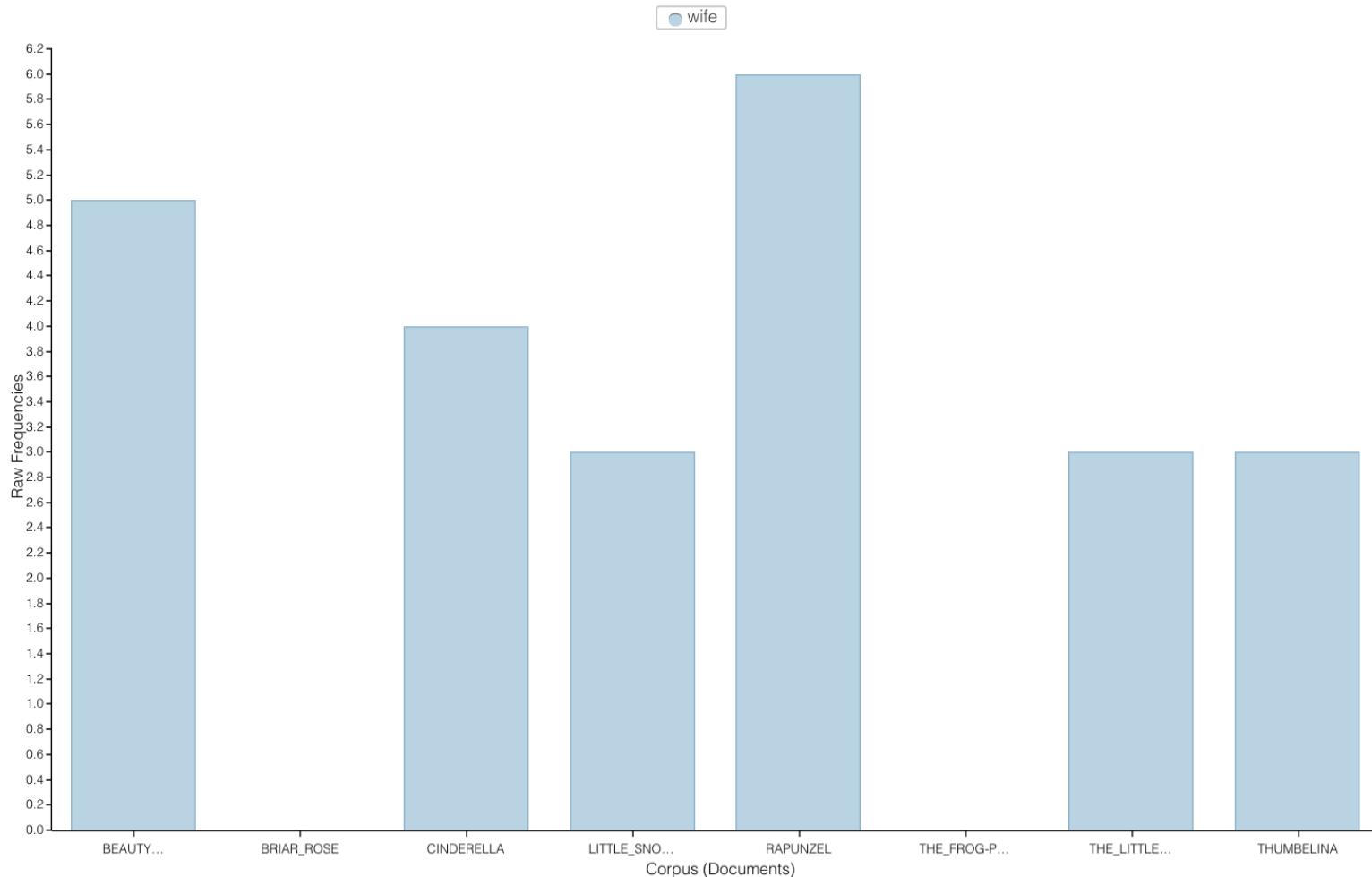
Mia and Erica are here to help!

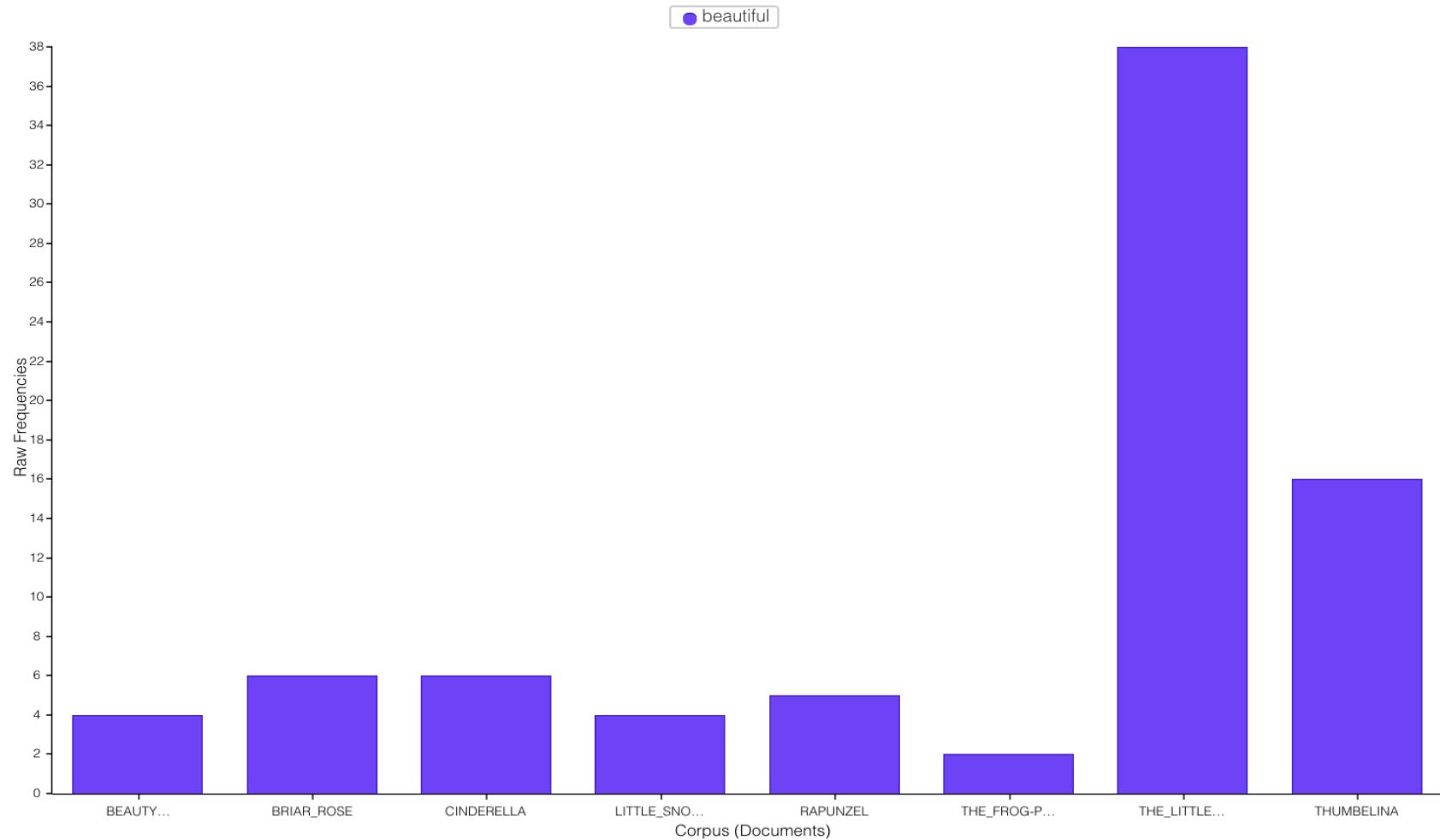
Raise a hand if you need help.



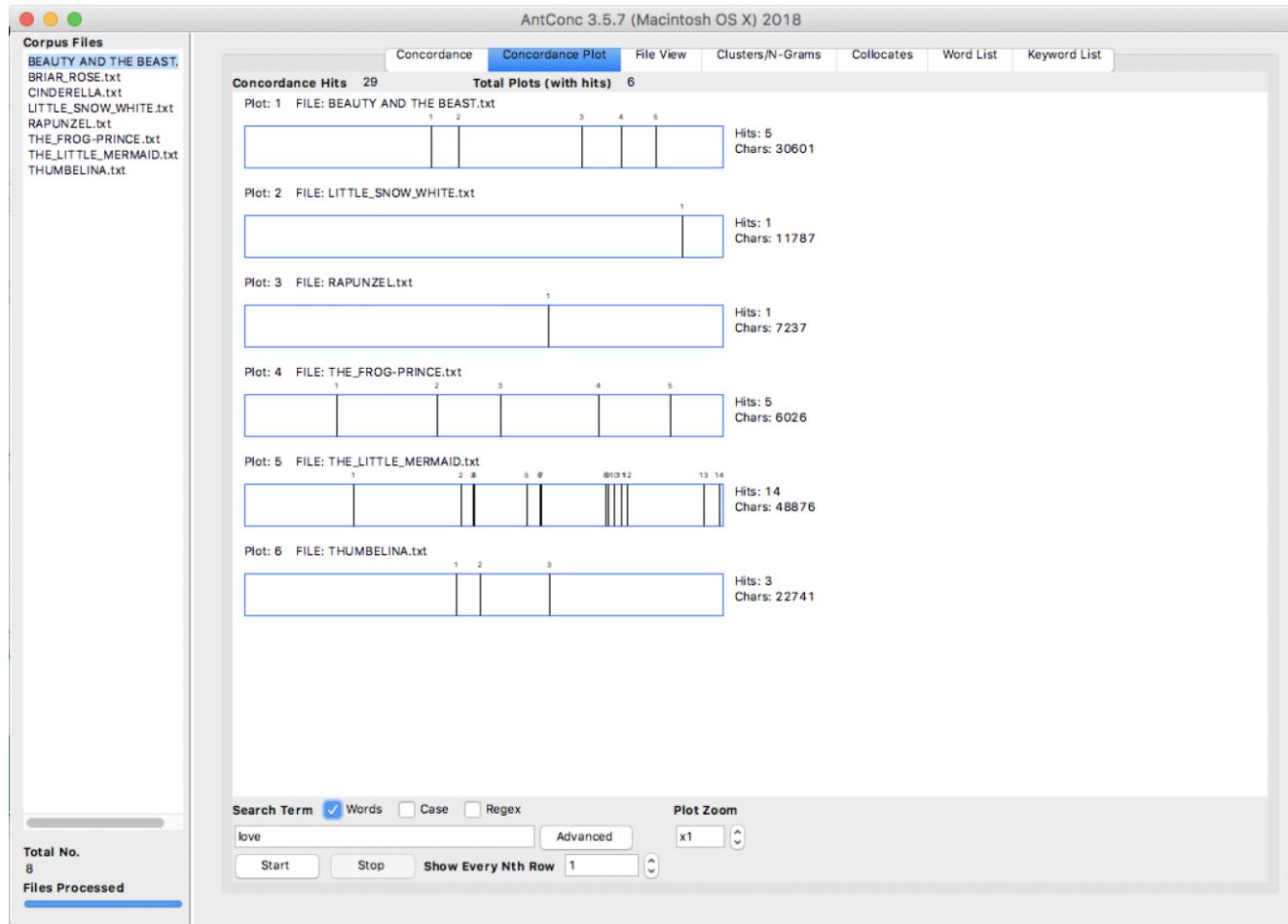








AntConc



AntConc 3.5.7 (Macintosh OS X) 2018

Corpus Files

BEAUTY AND THE BEAST.
BRIAR_ROSE.txt
CINDERELLA.txt
LITTLE_SNOW_WHITE.txt
RAPUNZEL.txt
THE_FROG-PRINCE.txt
THE_LITTLE_MERMAID.txt
THUMBELINA.txt

Concordance Concordance Plot File View Clusters/N-Grams Collocate Word List Keyword List

Concordance Hits 29

Hit	KWIC	File
1	bonly, and he loved her as he would love a little child, but it never came into his	THE_LITTLE_ME
2	only one in the world whom I could love; but you are like her, and you have almo	THE_LITTLE_ME
3	saving my father, and of proving my love for him." "No, my sister," said the three	BEAUTY AND T
4	ring her home as my bride. I cannot love her; she is not like the beautiful maide	THE_LITTLE_ME
5	ess dear, Open the door to thy true love here! And mind the words that thou a	THE_FROG-PRI
6	ess dear, Open the door to thy true love here! And mind the words that thou a	THE_FROG-PRI
7	ess dear, Open the door to thy true love here! And mind the words that thou a	THE_FROG-PRI
8	every day. I will take care of him, and love him, and give up my life for his sake."	THE_LITTLE_ME
9	has all these good qualities. I do not love him, but I respect him, and I feel both	BEAUTY AND T
10	wels, and fine clothes; but if you will love me, and let me live with you and eat	THE_FROG-PRI
11	nd handsome, she thought: 'He will love me more than old Dame Gothel does';	RAPUNZEL.txt
12	to the foam of the sea. "Do you not love me the best of them all?" the eyes of	THE_LITTLE_ME
13	, and Tiny nursed him with care and love. Neither the mole nor the field-mouse i	THUMBELINA.t
14	she obtain one unless she wins the love of a human being. On the power of and	THE_LITTLE_ME
15	ace again; and if you do not win the love of the prince, so that he is willing to	THE_LITTLE_ME
16	oy of his parents and deserves their love, our time of probation is shortened. Th	THE_LITTLE_ME
17	Tiny very sad to see it, she did so love the little birds; all the summer they had	THUMBELINA.t

Search Term Words Case Regex

love

Search Window Size

Start Stop Sort Show Every Nth Row

Kwic Sort Level 1 1R Level 2 2R Level 3 3R

Total No. 8
Files Processed

Online evaluation

go.ncsu.edu/libeval

Resources

- [Voyant Tools Documentation](#)
- [Voyant Tools Guide](#)
- [Digging through the Hillary Clinton's Archive with Voyant Tools](#)
- [Searching for the Victorians](#)
- [Ted Underwood's Where to Start with Text Mining](#)
- [AntConc Tutorial Videos](#)
- [The Programming Historian's Corpus Analysis with AntConc](#)