

# Introduction to Computational Text Analysis

---

INSH 2102: Bostonography  
Prof. Parr and Prof. Nelson  
Spring 2023

Juniper Johnson and Ana Abraham  
**Digital Integration Teaching Initiative (DITI)**



**Northeastern University**  
*NULab for Texts, Maps, and Networks*

*Feel free to ask questions at any point  
during the presentation!*

# Workshop Agenda

- Introduction to key terms and concepts in computational text analysis (What)
- Discussion on its applications and uses in research (Why)
- Demonstration of web-based text analysis tools (How)
  - Word Counter, Word Trees, Voyant, Lexos

Slides, handouts, and data available at:

<https://bit.ly/sp23-parr-textanalysis>



# What is Computational Text Analysis?



Northeastern University  
*NULab for Texts, Maps, and Networks*

*Feel free to ask questions at any point  
during the presentation!*

# Computational Text Analysis

Computational text analysis refers to the **array of methods used to “read” texts with a computer**. It is similar to statistical analysis, but the data is texts (words) instead of numbers.

Text analysis:

- Involves a computer drawing out patterns in a text, and a researcher interpreting those patterns.
- Includes methods such as word count frequency, keywords in context, computational modeling (with machine learning), and sentiment analysis.
- Is conducted using web-based tools or coding languages like Python and R



# Why Computational Text Analysis?

Computational text analysis can help us **analyze very large amounts of data, identify keywords, and discover patterns** in texts. Using text analysis, researchers may **find surprising results** that they would not have discovered from close reading or traditional methods alone.

For example: "[Gendered Language in Teacher Reviews](#)" by Ben Schmidt shows stark differences in the ways that male and female professors are reviewed on "Rate My Professor."



# Gendered Language



Go to [bit.ly/schmidt-gender](https://bit.ly/schmidt-gender) and try a few queries.

For example:

- Smart
- Ditz
- Unprofessional
- Nice
- What else did you try?
- How do you think Schmidt determined professor gender?



# Key Terms

- **Corpus (plural–corpora):** A collection of texts used for analysis and research purposes.
- **Stop words:** Words that appear frequently in a language, like pronouns, prepositions, and basic verbs. These are often removed for computational analysis. Some English stop words include: a, the, she, he, I, me, us, of, is, would, could, should, etc.
- **Word Count Frequency:** Counting the total times a word appears in a text/corpus or the percentage of how often it appears.
- **nGram:** A continuous sequence of  $n$  items in a text. A bigram (or 2 continuous words) could be ‘United States,’ while a trigram (3 words) could be ‘yes we can.’



# Corpus Building

## Questions to consider as you begin your research:

- What are my research questions and why am I creating a corpus?
- What am I asking my corpus to do?
- What text(s) should form my corpus to answer my research questions?
- How should I organize my corpus to streamline my research processes and save time?

For more on building a corpus, see [this handout](#).





# Our Corpus

For our corpus, we will work with a set of State of the Union addresses from 1990 to 2023.

You can download these files here: <https://bit.ly/insh2102-SOTU-corpus>

The easiest way to work with these files is to choose "Download all" and open them with a plain-text editor (TextEdit on Mac, Notepad on Windows). Mac users should be able to click on the zip file to open it; Windows users will need to right-click and choose "Extract all."



# Initial Corpus Analysis

**Open any one of the texts from the sample corpus:**

What can you observe about the text? How long is it? What kinds of language does it use? What kinds of analysis might you do with a text like this?

Scan through a few more: do they seem largely similar? What do you think might be different?



# Exploratory Tools: Word Counter and Word Trees



Northeastern University  
*NULab for Texts, Maps, and Networks*

*Feel free to ask questions at any point  
during the presentation!*

# Word Counter

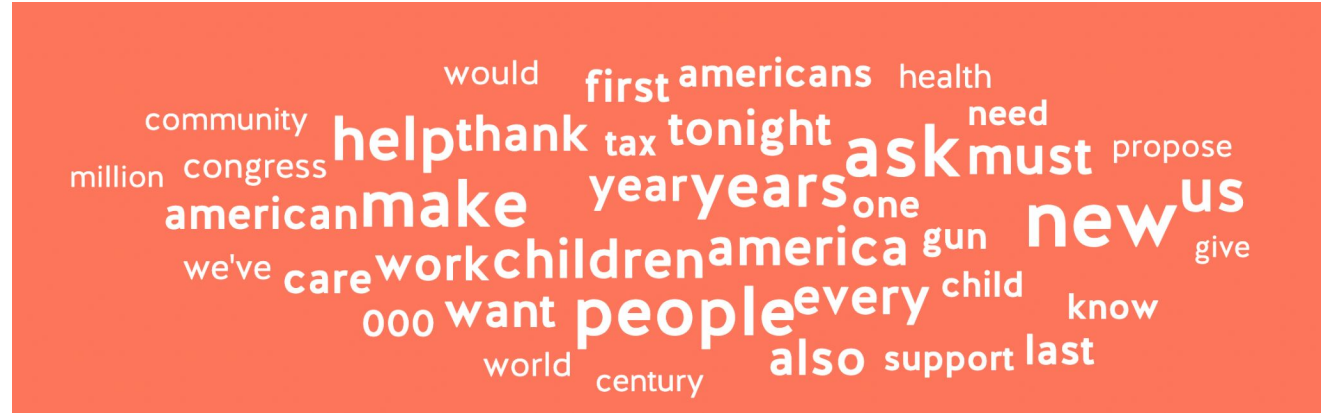
- <https://databasic.io/en/wordcounter/>
- A user-friendly **basic word counting tool**
- Allows you to count words, bigrams, and trigrams in plain text files and to download spreadsheets with your results
- The max file upload is 10MB
- The default is to lowercase all words and apply stopwords
- It can be run with and without stopwords and lowercasing



# Word Counter Example

This is a **word cloud**. It is helpful to get a sense of the **most used words** in a document.

Words used more often are bigger, and ones used less often are smaller.



## What seems significant in the most frequent terms from Clinton's 2000 SotU address?

Are there any things here that surprise you?

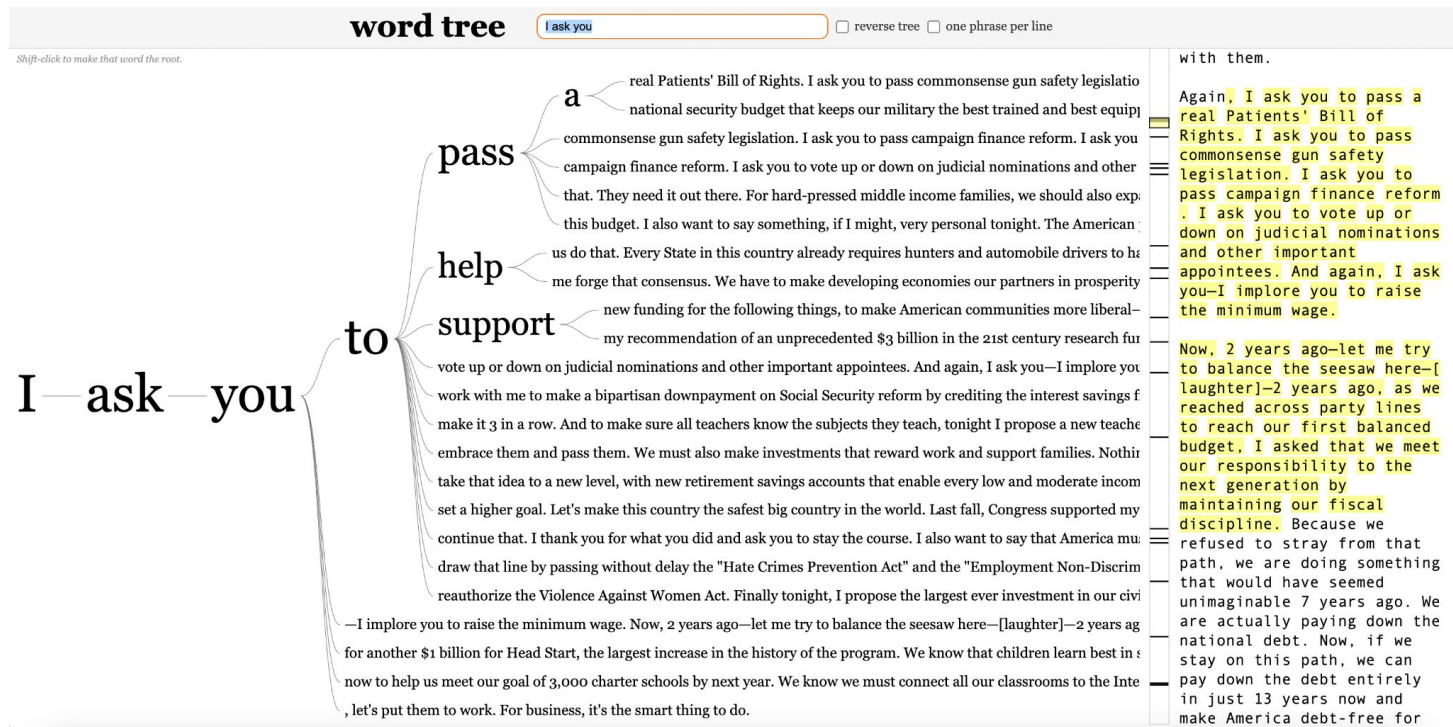


# Word Trees

- <https://www.jasondavies.com/wordtree/>
- A word tree **depicts multiple parallel sequences of words**
- This is a good way to see patterns in word usage, based on words that appear before and after a term or terms of interest.
- There are some restrictions in size with this tool: fewer than 1 million words should work



# Word Tree Example



# Your Turn!

Use the sample text or texts of your choice and begin practicing web-browser text analysis. **Explore Word Counter and Word Trees!**

## Discussion Prompts

- What limitations are you observing? What functionalities do you wish these tools might offer?
- Even with these limitations, how can you apply these simple tools in your research and exploration?





# Tools for corpus exploration: Voyant and Lexos



Northeastern University  
*NULab for Texts, Maps, and Networks*

*Feel free to ask questions at any point  
during the presentation!*

# Voyant

Voyant makes it possible to **perform analyses on one or multiple files in many ways**, including word counts, nGrams (n=number of words), word frequency distributions, word trends across documents, and concordances.

<https://voyant-tools.org/>



# VOYANT

see through your text

Click on Upload and navigate to the folder with the text documents you wish to analyze.

Alternatively, insert URLs or full text into the textbox.

Click here for help and advanced options

Add Texts

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

Open Upload Reveal





# Voyant: Changing displayed results

Hover on the right top corner of a pane and buttons will appear. Select the panes button and choose a new option from the dropdown menu. For example, we might want to try out the "Collocates" tool instead of the word cloud. Click on the ? to learn more about how the tool works.



	Terms	Links	Collocates	
	Term		Collocate	Count (context)
<input type="checkbox"/>	american		people	138
<input type="checkbox"/>	new		new	78
<input type="checkbox"/>	make		sure	77
<input type="checkbox"/>	new		jobs	71
<input type="checkbox"/>	years		ago	69
<input type="checkbox"/>	american		american	42
<input type="checkbox"/>	america		united	39
<input type="checkbox"/>	america		states	39
<input type="checkbox"/>	year		year	38
<input type="checkbox"/>	work		people	36
<input type="checkbox"/>	world		america	36
<input type="checkbox"/>	people		work	35



# Lexos

Lexos provides a step-by-step guide for text uploading, preparation, and analysis.

- **Upload:** upload your .txt file
- **Manage:** select the files you want to prepare and analyze
- **Prepare:** prepare your text for analysis
- **Visualize:** create visualizations of patterns across your corpus or in single texts
- **Analyze:** analyze your text

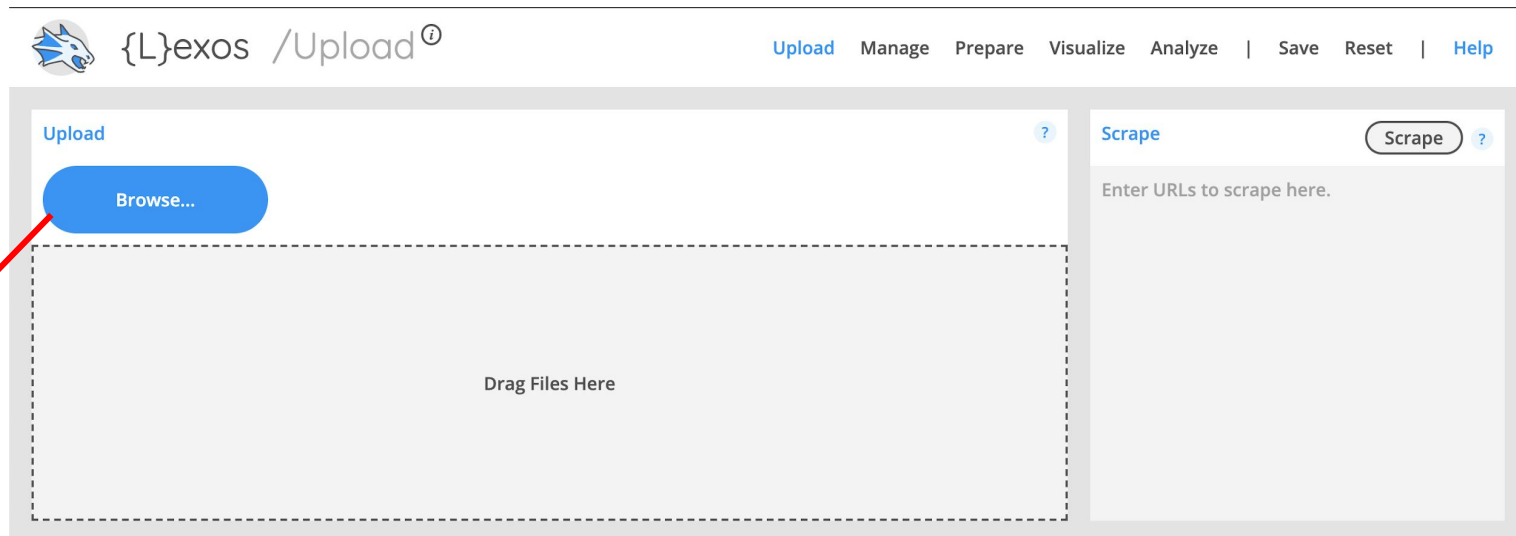
<http://lexos.wheatoncollege.edu/upload>



# Lexos: Upload


Click Browse and select your entire text (or drag file into the “Drag Files Here” area)

It can be easy to miss when the upload is done—click “Manage” to double check that the text file is there.







# Lexos: Manage

Make sure the document you want to use is selected (blue = selected, gray = not selected)

{L}exos /Manage<sup>i</sup>

Upload **Manage** Prepare Visualize Analyze | Save Reset | [Help](#)

Active	#	Document	Class	Source	Excerpt	Download ?
	3	2018-trump		2018-trump.txt	The President. Mr. Speaker, Mr. Vice President, Members of Congress, the First Lady of the United States, and my fellow America... ..er. And our Nation will forever be safe and strong and proud and mighty and free. Thank you. And God bless America. Goodnight.	
	4	2019-trump		2019-trump.txt	Madam Speaker, Mr. Vice President, Members of Congress, the First Lady of the United States — (applause) — and my fellow Americ... ..ong all the nations of the world. Thank you. God bless you. And God bless America. Thank you very much. Thank you. (Applause.)	
	5	2013-obama		2013-obama.txt	Please, everybody, have a seat. Mr. Speaker, Mr. Vice President, Members of Congress, fellow Americans: Fifty-one years ago, Jo... ..thors of the next great chapter of our American story. Thank you. God bless you, and God bless these United States of America.	
	6	2014-obama		2014-obama.txt	The President. Mr. Speaker, Mr. Vice President, Members of Congress, my fellow Americans: Today in America, a teacher spent ext... ..es cast toward tomorrow, I know it is within our reach. Believe it. God bless you, and God bless the United States of America.	





# Lexos: Prepare (scrub)

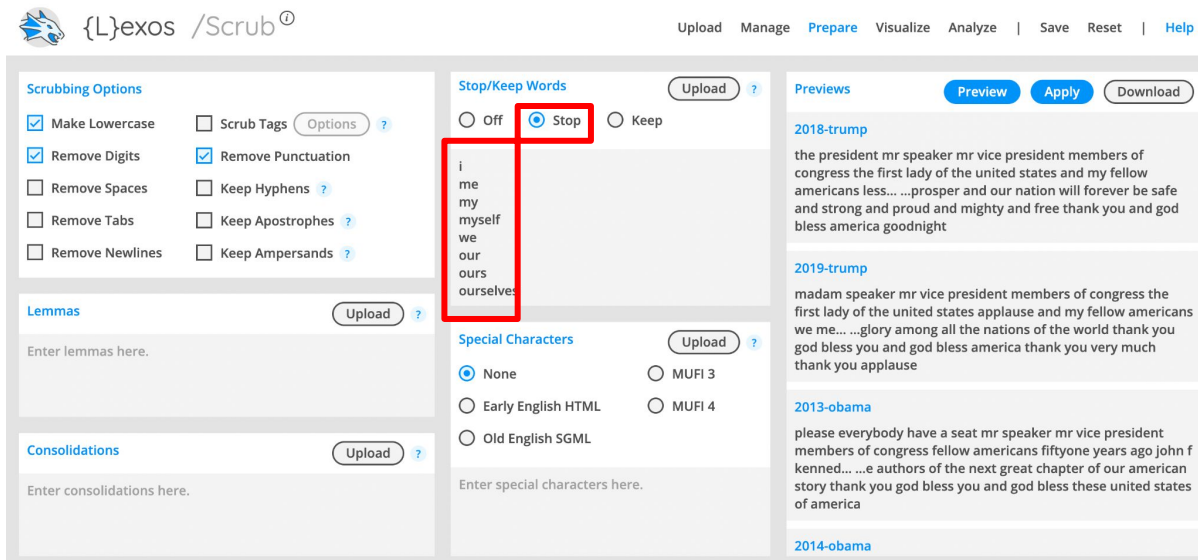
Lexos demonstrates some more advanced options you have for preparing your corpus. By “scrubbing,” you are transforming the texts in your corpus and making choices that will impact your results. Here are some possibilities:

- **Make Lowercase:** make all your letters lowercase. Even though you know “A” and “a” are the same letter, the computer treats these as two separate characters. Lowercasing removes this distinction.
- **Remove Punctuation:** remove punctuation, which may influence your results.
- **Stop/Keep Words:** remove a list of words. Usually these would be **stopwords**, or the most common words in a language (English: the, a she, her, it, him, they, etc).
- **Lemmas:** standardize to the *stem* of word. For example, you can stem all forms of talk: talking, talked, talks, etc. to “talk”



# Lexos: Removing Stopwords

Get a list of English stopwords here: <https://gist.github.com/sebleier/554280>. Copy and paste the stopwords (hit "raw", then select all and copy) into the “Stop/Keep Words” box then select “Stop”



The screenshot shows the Lexos Scrub interface. The top navigation bar includes links for Upload, Manage, Prepare, Visualize, Analyze, Save, Reset, and Help. The main interface is divided into several sections:

- Scrubbing Options:** Contains checkboxes for Make Lowercase, Remove Digits, Remove Spaces, Remove Tabs, Remove Newlines, Scrub Tags, Remove Punctuation, Keep Hyphens, Keep Apostrophes, and Keep Ampersands. Each option has a corresponding 'Options' link.
- Lemmas:** A section for entering lemmas with an 'Upload' button.
- Consolidations:** A section for entering consolidations with an 'Upload' button.
- Stop/Keep Words:** A section with radio buttons for Off, Stop, and Keep. The 'Stop' option is selected and highlighted with a red box. Below the radio buttons is a text area containing a list of stopwords: 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', and 'ourselves'. This text area is also highlighted with a red box.
- Special Characters:** A section with radio buttons for None, MUFI 3, MUFI 4, Early English HTML, and Old English SGML. The 'None' option is selected. There is an 'Upload' button and a text area for entering special characters.
- Previews:** A section with buttons for Preview, Apply, and Download. It displays three preview cards for the text '2018-trump', '2019-trump', and '2013-obama', showing the result of the scrubbing process.



# Lexos: Applying your Preparations

## BEFORE PREP

### 2013-obama

Please, everybody, have a seat. Mr. Speaker, Mr. Vice President, Members of Congress, fellow Americans: Fifty-one years ago, John F. Kennedy... the authors of the next great chapter of our American story. Thank you. God bless you, and God bless these United States of America.

### 2014-obama

The President. Mr. Speaker, Mr. Vice President, Members of Congress, my fellow Americans: Today in America, a teacher spent extra time... our eyes cast toward tomorrow, I know it is within our reach. Believe it. God bless you, and God bless the United States of America.

## AFTER PREP

### 2013-obama

please everybody have a seat mr speaker mr vice president members of congress fellow americans fiftyone years ago john f kenned... e authors of the next great chapter of our american story thank you god bless you and god bless these united states of america

### 2014-obama

the president mr speaker mr vice president members of congress my fellow americans today in america a teacher spent extra time... our eyes cast toward tomorrow i know it is within our reach believe it god bless you and god bless the united states of america

Once you have made decisions about your preparations, click “**Apply**” and wait a few minutes. Because the program is going through each document and completing all the processes you selected, it needs some time. Then, you will see the final results of your preparation! You can also **download** your new corpus.



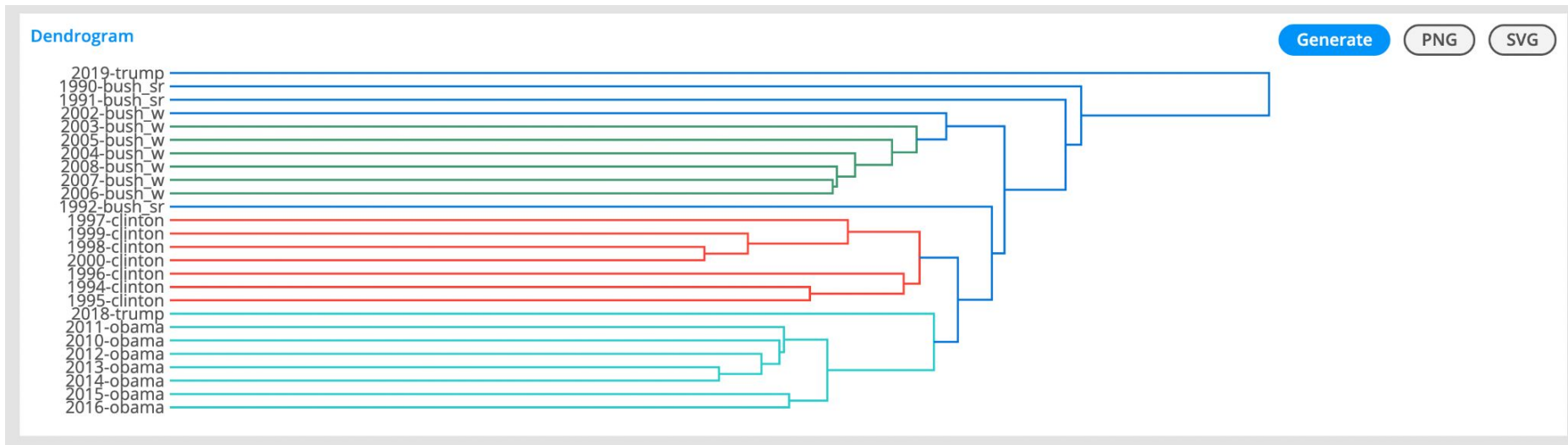
# Lexos: Analyze > Dendrogram

The dendrogram demonstrates similarity between the different documents. Dendrograms require at least two documents to compare. Dendrograms are able to show the hierarchy between objects. Dendrograms show:

- Similarities between texts
  - The greater the distance between texts, the less similar they are
  - The smaller the distance between texts, the more similar they are



# Lexos: Analyze > Dendrogram



# Lexos: Analyze > Top words



{L}exos /Top Words<sup>i</sup>

Upload Manage Prepare Visualize **Analyze** | Save Reset | [Help](#)

## Top Words

Generate

Download



### Document "2018-trump" Compared To The Corpus

cj	8.7532
ryan	8.4414
isis	8.0021
corey	7.9905
kenton	7.9905
preston	7.9905

### Document "2019-trump" Compared To The Corpus

applause	31.7392
usa	8.9133
elvin	8.6778
alice	8.2841
thank	8.1019
border	8.0326

### Document "2013-obama" Compared To The Corpus

desiline	6.5217
vote	6.4658
reduction	6.2286
preschool	6.0796
brian	5.6479
task	5.5521

### Document "2014-obama" Compared To The Corpus

cory	9.6023
workforce	5.6954
amanda	5.5435
easy	5.2962
irans	5.2681
equalitytv	5.0525



Northeastern University  
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point during the presentation!*

# Lexos: Save or Reset Your Results

Lexos allows you to **save** your results as a Lexos file. If you do this, you can re-upload the Lexos file any time to access your cleaned-up corpus as well as the different analyses you've done. You can also download modified text files from the “Manage” page, which you can use with other tools if you would like.

You can also save individual visualizations as images (PNGs).

Finally, if you want to start over, you can “Reset” your Lexos dashboard.



# Your Turn!

Use the sample text or texts of your choice and begin practicing web-browser text analysis. **Explore Lexos and Voyant!**

## Discussion Prompts

- What difference did you notice between Voyant and Lexos?
- Which tool do you prefer and why?
- How would you want to use these tools in this class and future?





# Learn more

Handouts (download these to use embedded links):

- [Building a corpus](#)
- More [links and resources](#) for text analysis
- [Lexos](#)
- [Voyant](#)
- [WordCounter](#)



# Thank you!

If you have any questions, contact us at [nulab.info@gmail.com](mailto:nulab.info@gmail.com)

**Sign up for our office hours at:** <https://calendly.com/diti-nu>

We'd love your feedback! Please fill out a short survey here:

<https://bit.ly/diti-feedback>

**Slides, handouts, and data available at:**

<https://bit.ly/sp23-parr-textanalysis>

