

Computational Text Analysis for Content Analysis

By Vaishali Kushwaha and Julianna Wessels
Digital Integration Teaching Initiative (DITI)

For ENGW 4710 Capstone Seminar
Mya Poe
Fall 2022



Northeastern University
NULab for Texts, Maps, and Networks

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

Workshop Agenda

- Introduction to definitions and key terms 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

Slides, handouts, and data available at:

<https://bit.ly/3de0GBh>



Workshop Objectives

- Understand best practices for collecting and storing textual data when performing basic computational text analysis
- Understand how web-based computational text analysis programs work, such as in their behind-the-scenes data preparation
- Understand how to interpret the results from your text analysis explorations



Workshop Outline

- Introduction
- Examples from Practice
- Text Preparation
- Exploratory Tools: Word Counter, Word Trees
 - Demo
- Voyant
 - Demo
- Lexos
 - Demo
- Conclusion



Introduction



Northeastern University
NULab for Texts, Maps, and Networks

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

Computational Text Analysis

Text analysis is making inferences based on textual data.

Computational text analysis refers to an array of methods that can be used to “read” texts with a computer. It is similar to statistical analysis, but the data are texts.

- It involves a computer drawing out patterns in a text, and a researcher interpreting those patterns.
- It includes methods such as word count frequency, nGrams, and sentiment analysis.



Why Computational Text Analysis?

Computational text analysis can help us **analyze very large amounts of data** and **discover patterns** in texts.

- From collections of texts, researchers can discover keywords that serve as a proxy for major trends in societies, cultures, and policies.
- Particular disciplines care deeply about the language used and how this language may reach intended audiences. Text analysis provides another method for approaching these discourses.
- Researchers may find surprising results that they would not have discovered from close reading or traditional methods alone.



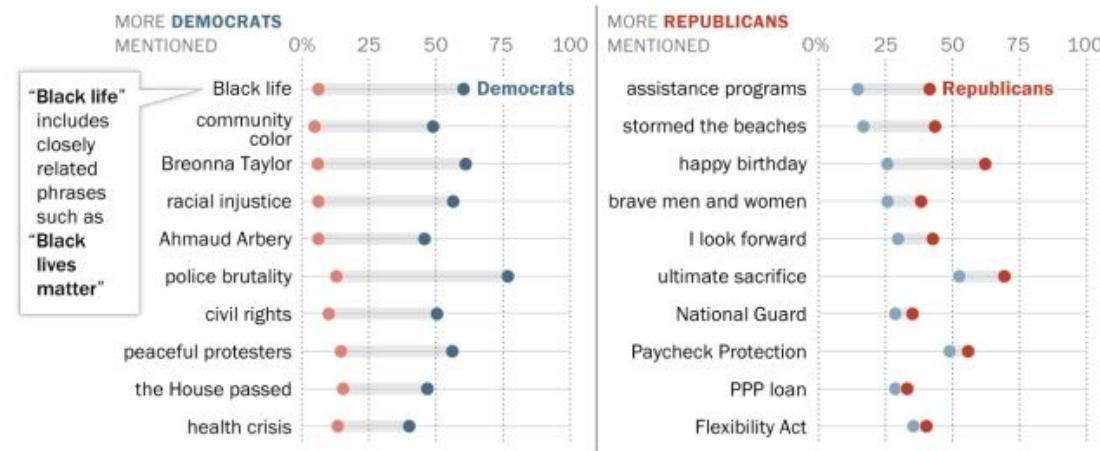
Content warning: police violence, racism

Posts mentioning 'Black lives matter' spiked on lawmakers' social media accounts after the death of George Floyd

- [Pew Research Center July 16, 2020 article](#)
- [Methodology](#)

In weeks following George Floyd killing, Democratic lawmakers' most distinctive language on social media focused on racial justice, police violence

Share of members in each party that mentioned ___ on Twitter or Facebook, May 25-June 14, 2020



Note: Chart shows the top 10 keywords based on how much more likely members of one party were to ever mention a keyword relative to the other party. Terms are displayed in their standardized form (e.g., "Black life" instead of "Black lives") and have been edited slightly in some cases for readability (e.g., "the House passed" instead of "house passed"). Keyword analysis was not case-sensitive. Words from retweets are included in this analysis even if the member who retweeted them did not create the original tweet.

Source: Pew Research Center analysis of congressional social media data from the Twitter API, Facebook Graph API and CrowdTangle, May 25-June 14, 2020.

PEW RESEARCH CENTER



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.’
- **Sentiment Analysis:** Measuring the sentiment of a text based on a scale such as negative/positive or happy/sad. Each word has a particular weight to determine where on the scale it falls, and these weights are calculated to determine a text’s overall sentiment.



Text Preparation



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?



Preparing Your Text

1. Choose the texts or text selections that you would like to include.
2. Create a folder on your computer or cloud storage where you will store your corpus. Give it a clearly descriptive name, without spaces or special characters.
3. Copy and paste the text into a **plain text editor** (on Macs: Text Edit; on Windows: Notepad)
 - a. Mac users, you may need to make your Text Edit into a ‘plain text’. Open Text Edit, go to Preferences, and make sure “plain text” is selected
4. Save the text as a plain text file (with a .txt extension). Always make sure to name your files so you know what is in them!
5. Repeat steps above for each text in the corpus.



Our Text

Our text is plain text (.txt file) of

Filippo Tommaso Marinetti, “[The Joy of Mechanical Force](#)” /
“[Futuristic Manifesto](#)” ([first edition](#))

McKenzie Wark, “[A Hacker Manifesto](#)”



Sample Corpus

The following .txt files are available on: <https://bit.ly/3de0GBh>



Exploratory Tool: Word Counter



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**
- It 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
- Default is lowercase all words and apply stopwords
- It can be run with and without stopwords



Word Counter Examples

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.



Word Counter Examples

TOP WORDS 	
Word	Frequency
class	98
information	82
property	59
production	40
form	39
politics	39
hacker	37
new	32
hack	31
free	31

Shows the top words in the text.

Stopwords aren't removed for the bigrams and trigrams because they need context.

BIGRAMS 	
bigram 	Frequency
of the	81
the hacker	33
to the	31
as a	31
is the	30
of information	25
in the	24
hacker class	22
the hack	20
of a	19

TRIGRAMS 	
trigram 	Frequency
the hacker class	22
the vectoralist class	13
the production of	11
of the hack	10
the possibility of	9
form of property	9
as a class	8
the means of	8
the form of	8
the politics of	8

The top two trigrams 'the hacker class' and 'the vectoralist class' both contain the words 'the' and 'class'. 'The' is a stopword, and 'class' is the dominant word in this text!



Exploratory Tool: Word Trees



Northeastern University
NULab for Texts, Maps, and Networks

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

Word Trees

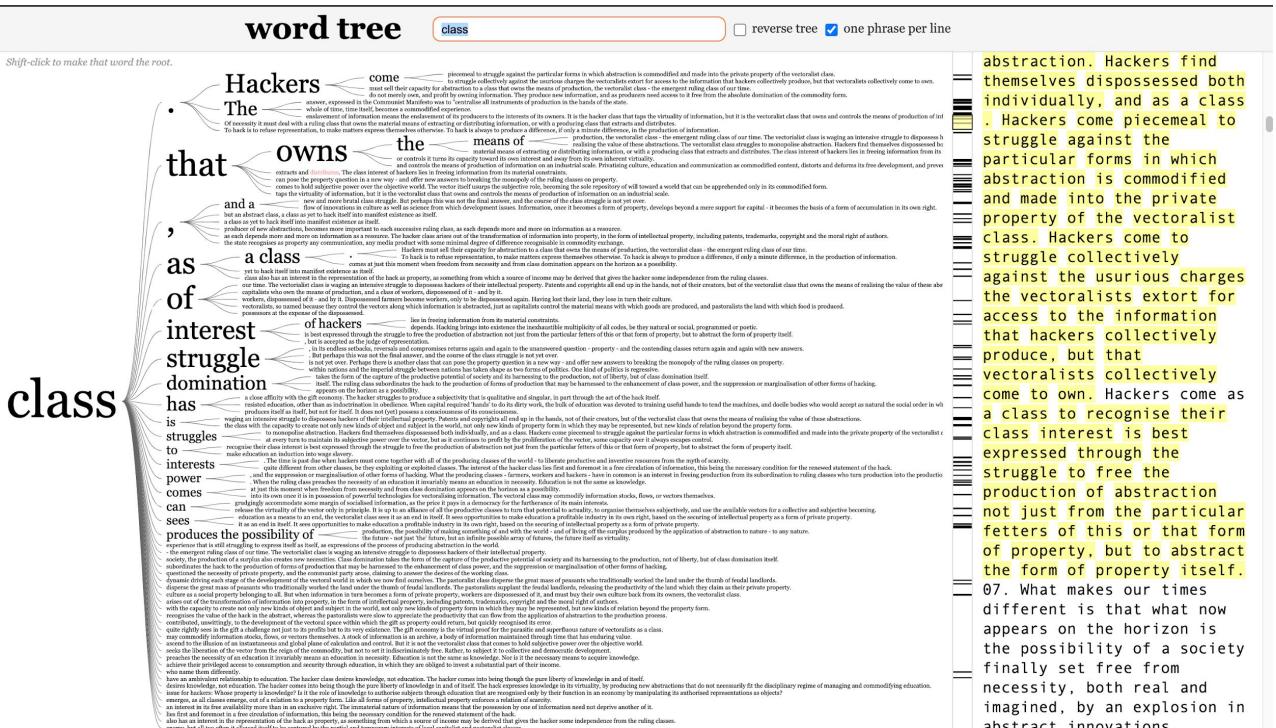
- <https://www.jasondavies.com/wordtree/>
- A word tree **depicts multiple parallel sequences of words**
- This is a good way to see patterns in word usage, based on words that appear before and after a term or terms of interest.
- There are some restrictions in size: fewer than 1 million words should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



Word Tree Example

Reflects the focus of the manifesto on class interest, class struggle, class domination, etc.

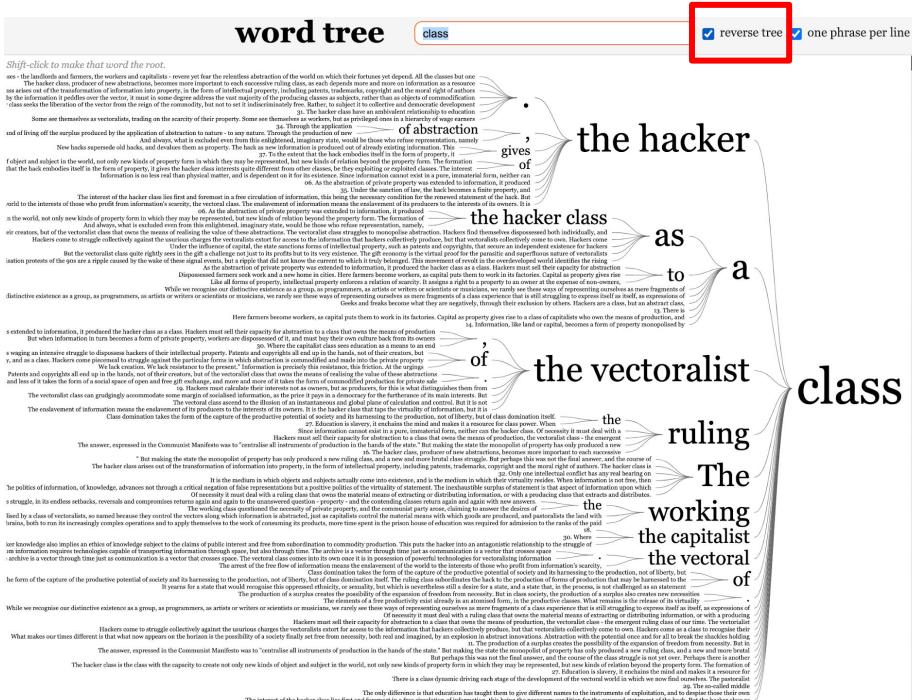
Notice the word
'class' seems to be
often at the end of a
sentence, followed by
period.



Word Tree: Reverse Trees

It is worth reversing the tree to see the words that often precede it. To do this click “reverse tree” next to the search bar.

Here the hacker, the vectoralist, ruling, working etc. are the dominant words preceding the word 'class.



fortunes of states and armies, companies and communities depend on it. All contending classes - the landlords and farmers, the workers and capitalists - revere yet fear the relentless abstraction of the world on which their fortunes yet depend. All the classes but one. The hacker class. Q2. Whatever code we hack, be it programming language, poetic language, math or music, curves or colourings, we create the possibility of new things entering the world. Not always great things, or even good things, but new things. In art, in science, in philosophy and culture, in any production of knowledge where data can be gathered, where information can be extracted from it, and where in that information new possibilities for the world are produced, there are hackers hacking the new out of the old. While hackers create these new worlds, they do not possess them. That which we create is mortgaged to others, and to the

Voyant



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. It also makes nice visualizations!

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



VOYANT

see through your text

Add Texts

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

?

Reveal

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

Alternatively, insert URLs or full text into the textbox.

Click here for help and advanced options



Northeastern University
NULab for Texts, Maps, and Networks

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

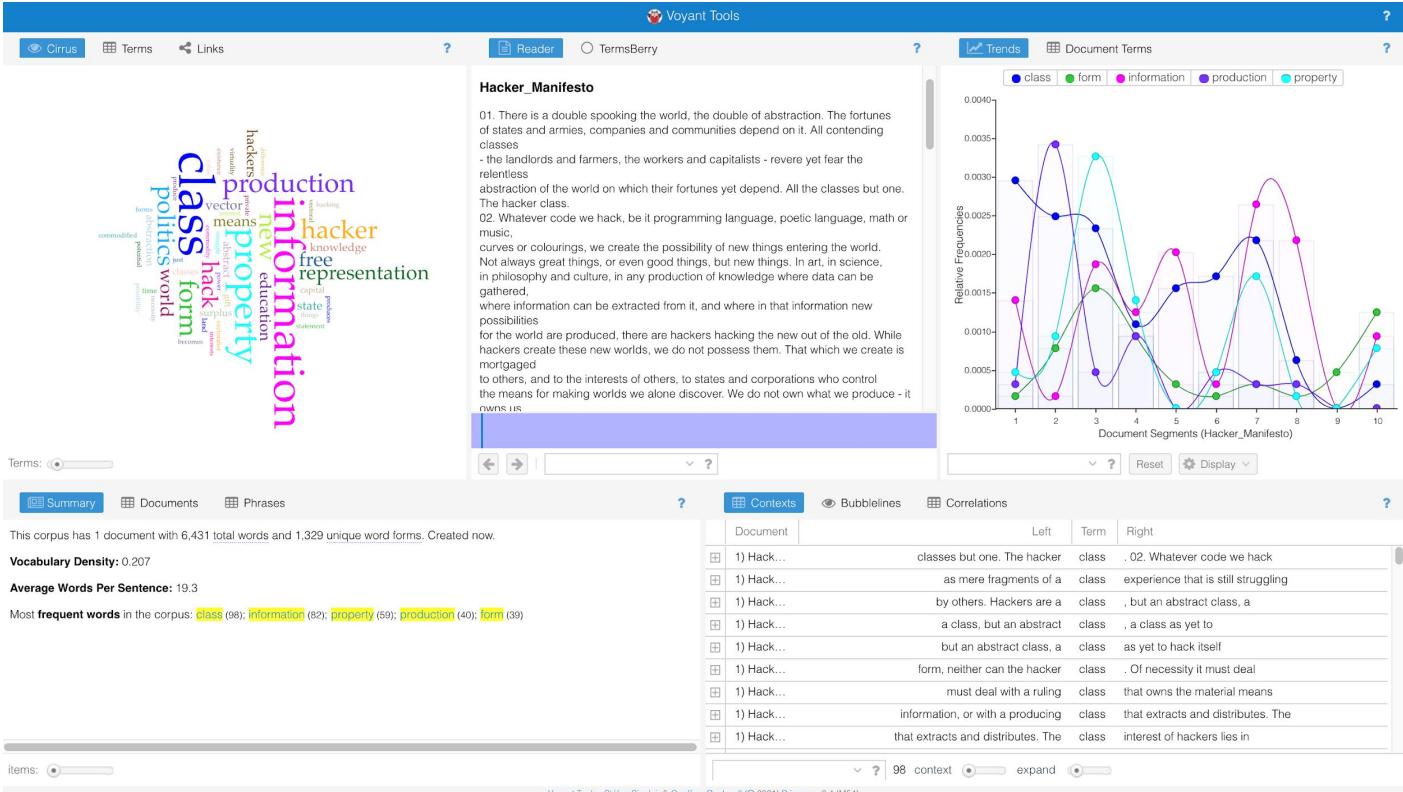
Voyant: Understanding the Dashboard

Results:

From Climate Ready Boston you can see the default results page with multiple panes:

- A word cloud
- Reader section
- Trends
- Document Summary
- Word Contexts

These boxes can all be changed!



Voyant: Contexts (concordances)

Contexts, or concordances, show the different contexts around particular search terms. For example, you can see all the times the word “class” appears in the text and the contexts in which it appears.

Contexts Bubblelines Correlations

Document	Left	Term	Right
1) Hacker_Man...	classes but one. The hacker	class	. 02. Whatever code we hack
1) Hacker_Man...	as mere fragments of a	class	experience that is still struggling
1) Hacker_Man...	by others. Hackers are a	class	, but an abstract class, a
1) Hacker_Man...	a class, but an abstract	class	, a class as yet to
1) Hacker_Man...	but an abstract class, a	class	as yet to hack itself
1) Hacker_Man...	form, neither can the hacker	class	. Of necessity it must deal
1) Hacker_Man...	must deal with a ruling	class	that owns the material means
1) Hacker_Man...	information, or with a producing	class	that extracts and distributes. The
1) Hacker_Man...	that extracts and distributes. The	class	interest of hackers lies in

▼ ? 98 context expand



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

The figure displays two word clouds side-by-side, connected by a red dashed arrow pointing from left to right.

Left Word Cloud: This cloud is centered around the word "class". Key terms include "new", "information", "property", "free", "form", "means", "politics", "production", "representation", "hack", "state", "education", "abstraction", "knowledge", "world", "vector", "vectorialist", "abstract", "time", "space", "difference", "becomes", "possibility", "existence", and "interest".

Right Word Cloud: This cloud is centered around the word "hacker". Key terms include "hacker", "abstract", "ruling", "producing", "vectoralist", "working", "another", "the", "vectoral", "as", "that", and "class".

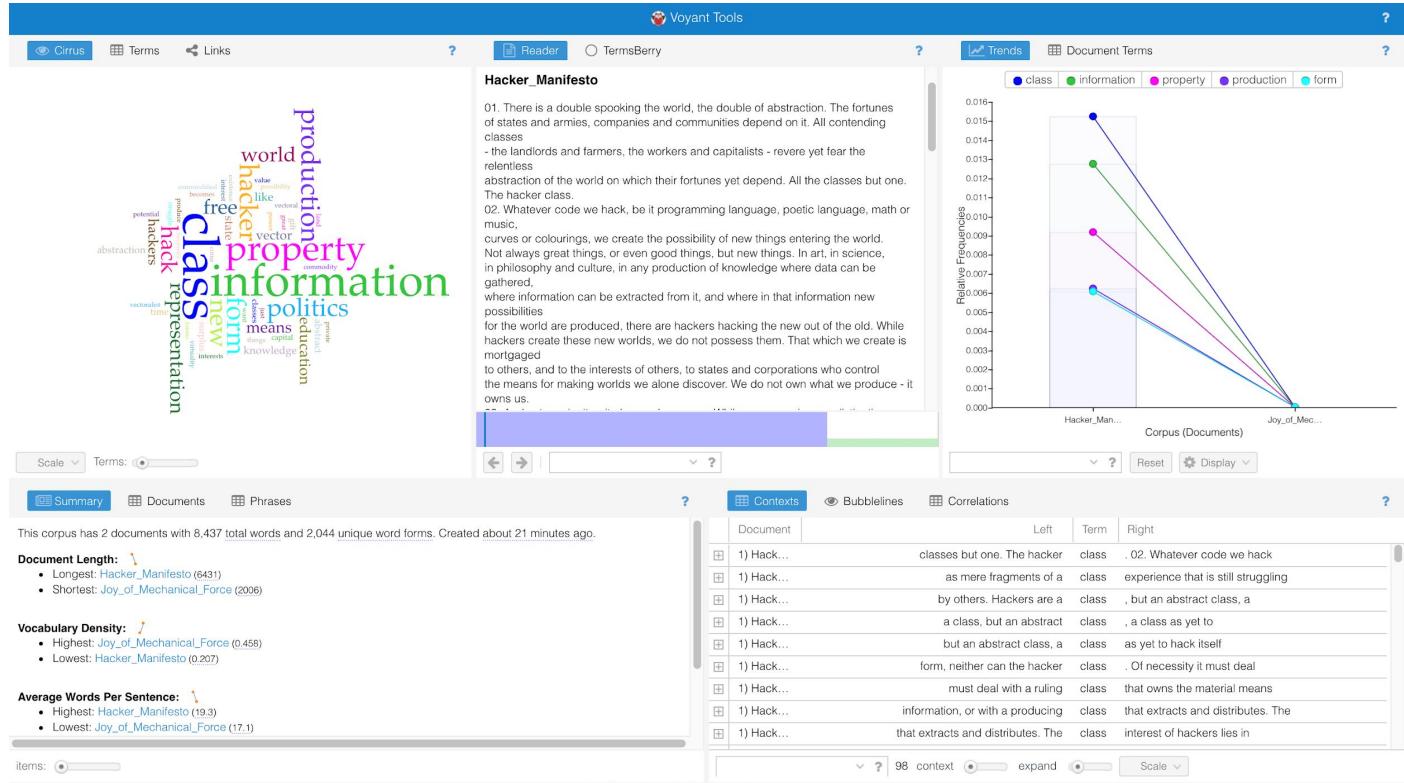
For our new pane option, we have chosen the WordTree visualization from the ‘visualization tools’ dropdown sub-menu. You can select the number of “branches” by dragging the scroll button at the bottom.



Voyant: Corpus Dashboard

Results page of the corpus containing climate reports of 5 cities.

- A word cloud: combining all texts
- Reader section: scroll down all texts
- Trends: relative frequency of terms across text - good for comparison
- Document Summary- good for comparison
- Word Contexts: separate for all texts



Lexos



Northeastern University
NULab for Texts, Maps, and Networks

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

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>



Northeastern University
NULab for Texts, Maps, and Networks

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

Lexos: Upload

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

You will not get a super visible notification when the upload is done - click “Manage” to double check that the text file is there.

The screenshot shows the Lexos interface with the 'Upload' tab selected. On the left, there's a 'Browse...' button and a dashed rectangular area labeled 'Drag Files Here'. A red arrow points from the text in the orange box to the 'Browse...' button. On the right, there's a 'Scrape' section with a 'Scrape' button and a placeholder text 'Enter URLs to scrape here.'



Lexos: Manage

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

{L}exos /Manage⁽ⁱ⁾

Upload **Manage** Prepare Visualize Analyze | Save Reset | Help

Active	#	Document	Class	Source	Excerpt	Download	?
<input checked="" type="radio"/>	1	Hacker_Manifesto		Hacker_Manifesto.txt	Â 01. There is a double spooking the world, the double of abstraction. The fortunes of states and armies, companies and commun... ...s to permeate existing states with a new state of existence, spreading the seeds of an alternative practice of everyday life.		
<input type="radio"/>	2	Joy_of_Mechanical_Force		Joy_of_Mechanical_Force.txt	The Joy of Mechanical Force by F. T. Marinetti ("The Foundation of Futurism" ["Manifesto of Futurism," 1909], translated from t... ...repeating these infamous words! Rather, look up! Up on the crest of the world, once more we hurl our challenge to the stars!		



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> (there is also a copy on the GitHub page). Copy and paste the stopwords (or upload the .txt file) into the “Stop/Keep Words” box then select “Stop”

The screenshot shows the Lexos Scrub interface. In the center, there is a 'Stop/Keep Words' section with three radio button options: 'Off' (unchecked), 'Stop' (checked), and 'Keep' (unchecked). Below this, a list of words is displayed: 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', and 'you'. To the right of this list are two small icons: a yellow face with a frown and a green circle with a checkmark. A red box highlights the 'Stop' radio button and the list of words. The top navigation bar includes 'Upload', 'Manage', 'Prepare' (which is blue, indicating it's selected), 'Visualize', 'Analyze', and buttons for 'Save', 'Reset', and 'Help'. On the left, there are sections for 'Scrubbing Options' (checkboxes for 'Make Lowercase', 'Remove Digits', etc.) and 'Lemmas' (an 'Upload' button and a text input field 'Enter lemmas here.'). On the right, there are 'Previews' for 'Hacker_Manifesto' and 'Joy_of_Mechanical_Force', each with a 'Preview' button, an 'Apply' button, and a 'Download' button.



Lexos: Applying your Preparations

BEFORE PREP

Previews Preview **Apply** Download

Hacker_Manifesto

Â 01. There is a double spooking the world, the double of abstraction. The fortunes of states and armies, companies and commun....s to permeate existing states with a new state of existence, spreading the seeds of an alternative practice of everyday life.

Joy_of_Mechanical_Force

The Joy of Mechanical Force by F. T. Marinetti ("The Foundation of Futurism" ["Manifesto of Futurism," 1909], translated from t... ...repeating these infamous words! Rather, look up! Up on the crest of the world, once more we hurl our challenge to the stars!

AFTER PREP

Previews Preview **Apply** Download

Hacker_Manifesto

â double spooking world double abstraction fortunes states armies companies communities depend contending classes landlords f... ...ing coalition interests seeks permeate existing states new state existence spreading seeds alternative practice everyday life

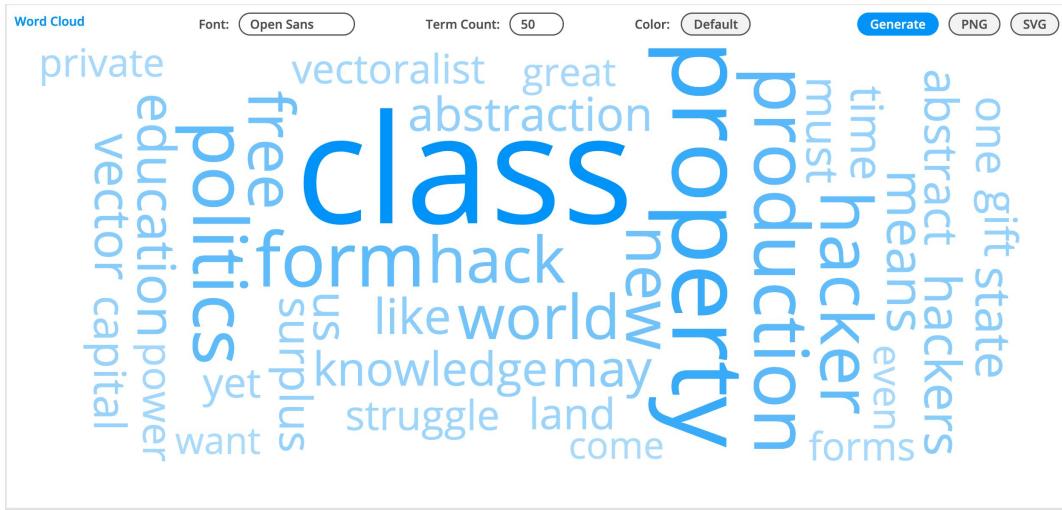
Joy_of_Mechanical_Force

joy mechanical force f marinetti foundation futurism manifesto futurism translated french eugen weber reprinted permission dod... ...on forebears perhaps let matter dont want listen beware repeating infamous words rather look crest world hurl challenge stars

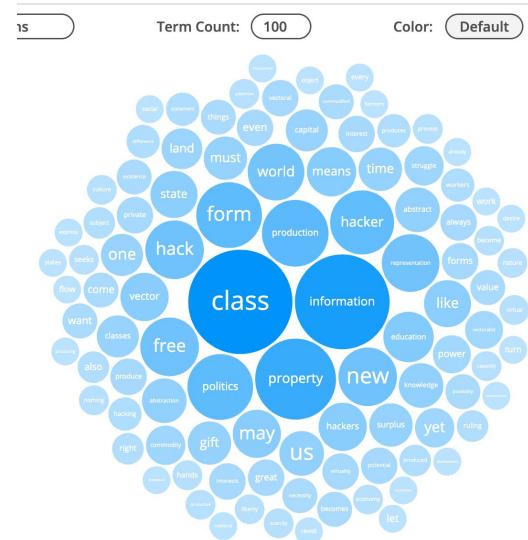
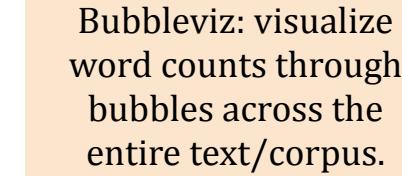
Once you have made decisions about your preparations, click “**Apply**” and wait a few minutes. Because the program is going through each document and completing all the processes you selected, it needs some time. Then, you will see the final results of your preparation! You can also **download** your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



Lexos: Visualize



Word Cloud: visualize a wordcloud across the entire text/corpus.



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



Lexos: Visualize > Multicloud



Northeastern University

NULab for Texts, Maps, and Networks

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

Lexos: Rolling Window

Rolling windows allow you to look at word trends across one document. To use a rolling window:

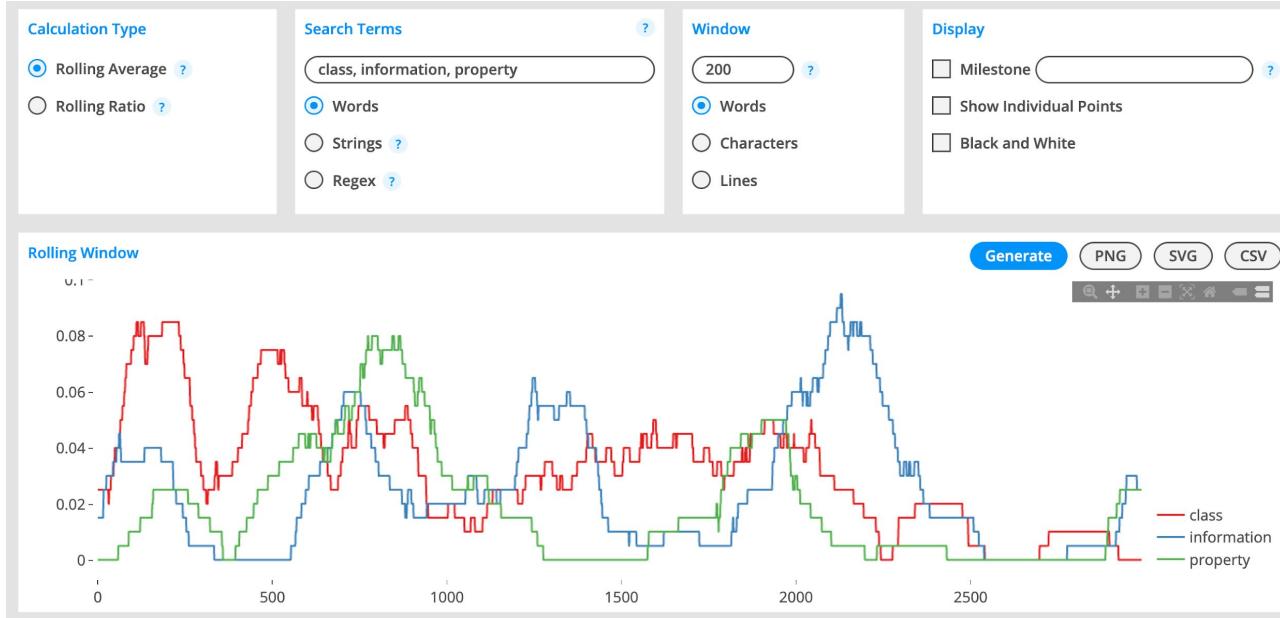
1. Go to “Visualize-> Rolling Window” and type in a search term you want to visualize. You can also search multiple terms by clicking “String” and separating words with a comma (heat, health, flood, storm)
2. Choose a Window size (the number of words each “window” contains). For shorter documents, it’s good to have a number like 300/500. For larger documents, you may want to make your window larger. Play around with the window size until you get a visualization that makes sense.
3. Click “Generate”

Calculation Type <input checked="" type="radio"/> Rolling Average ? <input type="radio"/> Rolling Ratio ?	Search Terms class, information, property <input checked="" type="radio"/> Words <input type="radio"/> Strings ? <input type="radio"/> Regex ?	Window 200 <input checked="" type="radio"/> Words <input type="radio"/> Characters <input type="radio"/> Lines	Display <input type="checkbox"/> Milestone ? <input type="checkbox"/> Show Individual Points <input type="checkbox"/> Black and White
--	---	---	---



Lexos: Rolling Window Results

Using *Hacker Manifesto*, and searching for the words ‘class’, ‘information’ and ‘property’ with a window of 200 (short document), we can get an idea of how these terms work together in the manifesto.



Lexos: Analyze > Dendrogram

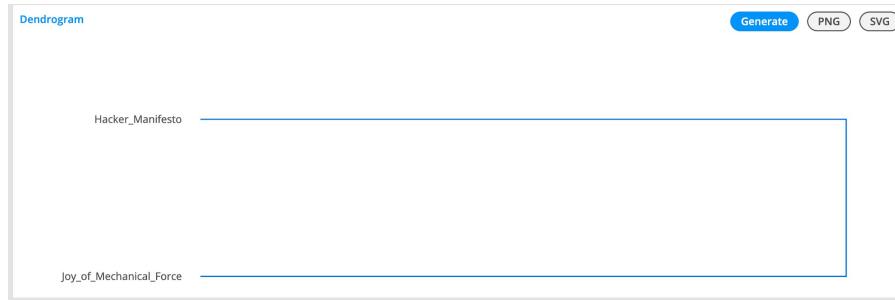
The dendrogram demonstrates similarity between the different documents. Dendograms require at least two documents to compare. Dendograms are able to show the hierarchy between objects. Dendograms show:

- 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: Dendrogram

The dendrogram demonstrates similarity between the different documents.



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.



Conclusion



Northeastern University
NULab for Texts, Maps, and Networks

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

Your Turn!

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

Discussion Prompts

- What do you find challenging or exciting about these tools?
- What interesting or surprising results came up?
- How might you interpret those results based on what you know about your field?



Thank you!

If you have any questions, contact DITI at nulab.info@gmail.com

Developed by:

Vaishali Kushwaha
DITI Teaching
Fellow

Julianna Wessels
NULab
Co-Coordinator

Milan Skobic
DITI Assistant
Director

Garrett Morrow
DITI Research Fellow

Slides, handouts, and data available at <https://bit.ly/3de0GBh>
Schedule an appointment with us! <https://calendly.com/diti-nu>



Northeastern University
NULab for Texts, Maps, and Networks

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