Mind Map Automation: Using Natural Language Processing to Graphically Represent
a Portion of a U.S. History Textbook


by
Sophie McIntyre


A thesis presented to the Honors College of Middle Tennessee State
University in partial fulfillment of the requirements for graduation
from the University Honors College

Spring 2021


Thesis Committee:


Dr. Salvador Barbosa, Thesis Director

Dr. Philip Phillips, Thesis Committee Chair

Mind Map Automation: Using Natural Language Processing to Graphically Represent
a Portion of a U.S. History Textbook

by Sophie McIntyre

APPROVED:

_____
Dr. Salvador Barbosa, Thesis Director
Assistant Professor, Computer Science

_____
Dr. Philip Phillips, Thesis Committee Chair
Associate Dean, University Honors College

**Acknowledgments**

I want to thank Dr. Barbosa for giving me the opportunity to work on a project that means a lot to me and for supporting me throughout this entire process. I couldn't have gotten to this point without his guidance. Additionally, I want to thank Dan for always listening to me whenever I needed to talk through a problem and for all the late nights where he stayed up with me while I worked. Finally, I want to let my parents, Wynne and Isabella know how much I appreciate their love and support.

**Abstract**

A mind map is a hierarchical representation of ideas that can be branched back to one centralized theme. Mind mapping can be applied to many disciplines and provides numerous cognitive benefits. One problem, however, is that the process of creating a mind map can become tedious and time consuming. This is especially problematic for students who have busy schedules and need to optimize their studying time. Therefore, this thesis aims to solve the challenge by developing an automated mind mapping system to analyze sections of a U.S. history textbook. The system uses natural language processing techniques to organize the text so that it can be graphically displayed for users. The primary goal of this project is to give students a supplementary tool that they can use to further their studies, and a future objective of this research is to expand the scope to include subjects outside of U.S. history.

**Table of Contents**

# List of Figures

## List of Terms

**Basic Ordering Ideas (BOIs)**: the first-level branches or chapter headings of a mind map (Buzan)

**JavaScript Object Notation (JSON)**: A data interchange and file format that was originally based on features found in the JavaScript programming language

**Library (programming)**: a compilation of predefined functions in which a computer program can use without having to explicitly state the function it calls

**Mind Mapping**: A method for displaying concepts that are derived from one central theme in a graphical format

**Natural Language Processing (NLP)**: a field of research that combines computer science with linguistics to achieve the goal of human language comprehension by machines

**Python**: high-level programming language

**Radiant Thinking**: thought process that starts with one centralized idea and naturally expands to include related concepts (Spencer 291)

**Repository**

All programs, texts and graphs can be found on GitHub through the following repository:

https://github.com/sdmac101/Mind-Mapping

**Chapter 1: Introduction**

Mind maps are organizational tools that have been gaining popularity over the past 50 years. This is in part due to the broad applicability of the mind mapping method within multiple disciplines, like education and business, as well as the positive effects it has on cognitive abilities. Although utilizing mind maps in different areas of one's life can be beneficial, the steps to build a proper map are repetitive and, depending on the subject matter, time consuming. With the continuous advancement of technology, when a problem like this is encountered, the search for a solution often begins by utilizing the capabilities of a computer to perform tedious jobs. Therefore, this project seeks to use modern computational technologies and applications to automate the creation of mind map inspired graphs. Since computers are unable to process information in the same way humans do, special methods must be used to extract information within textual data. The primary tools used within this project are natural language processing techniques, such as text summarization, coreference resolution and parsing, along with existing data visualization tools. The text being examined comes from an opensource U.S. history textbook, and the result of this project is a visual layout of all major concepts in a specified period of history in addition to the relationships between those concepts.

A mind map is a useful tool in many different fields, but it particularly excels in education. Mind maps are another form of note-taking, which is an integral part in most learning environments. However, many students, primarily ones in college, might stay away from mind mapping because it takes too much time or because they feel more comfortable with standard note-taking, despite the intellectual benefits mind maps provide. By simplifying the creation of mind maps, students will be able to gain the

advantages they offer while maintaining a busy schedule. Teachers can also use an application like this to help determine which concepts should be taken from a textbook and incorporated into the classroom. Through this project, a product is developed that can help save time for individuals in education while still providing a comprehensive overview of various topics. It is important to note that this project only focuses on examining history textbooks, but a future goal would be to make it applicable to other subjects.

Currently, digital mind map visualization tools, such as iMindMap, are available for anyone to use. However, the user is still required to manually enter the content of the map, and, for many people who prefer handwriting to typing, creating a digital version of their mind map may be even more tedious than drawing it by hand. Although there have been other studies regarding automated mind maps, they often use outdated techniques and do not specialize in educational materials. This project uses modern natural language processing methods and up-to-date visualization tools that require minimal input from the user.

## Chapter 2: Background and Related Works

### 2.1 Mind Maps

After studying the psychological effect that note-taking has on memory, Tony Buzan developed the mind mapping method in the 1970s (Buzan and Griffiths). Buzan first noticed problems with standard note-taking while in high school when he began to have difficulty recalling lessons he had learned in class. He began to underline and box important ideas, which helped a small amount, but the problem resurfaced when he was in college (Buzan 4). Buzan then found inspiration in the ancient Greeks, who had created elaborate mnemonic devices through the use of imagination and association. While continuing his research, he studied the psychology of thinking and discovered that the two primary contributors to memory were the same concepts used by the Greeks. This prompted him to create a thinking tool that would mimic the way a human brain processes and recalls information. Once finished, he presented his findings as the mind mapping method in the late 1960s. Shortly after this, BBC reached out to Buzan to do a 10-episode series where he defined what a mind map is and explained how a viewer could apply the method to their life. Buzan also wrote an accompanying book to the series titled *Use Your Head* which helped mind mapping gain significant popularity (Buzan 5). Buzan, who passed away in 2019, made significant contributions to the field of memory and literacy that have benefitted many other fields of study.

Mind maps are graphical representations that merge creativity with logical reasoning in order to promote cognitive thinking and improve memory (Mento 390). The structure of a mind map is based on Radiant Thinking, which dictates that every detail in the map "radiates" from a central point ("Mind Mapping or Concept Mapping"). This

means that concepts are organized by hierarchy in relation to the main idea. For instance, if the main subject of a mind map is food, it would be a reasonable assumption to believe that many would write the following food groups, meat, dairy, grains, vegetables, and fruit, as secondary thoughts that radiate from the general concept of food. By utilizing Radiant Thinking and imaginative features, mind mapping takes a "whole-brain" approach. That is, it incorporates the creativity of the right side of the brain with the rationality of the left side of the brain, resulting in an improvement to association and subsequently memory (Buzan 244).

The steps involved in producing a mind map are relatively simple, and the design process is meant to mimic the structure of our brain (Buzan and Griffiths). The first step in the design is to identify an image or keyword that will be located in the center of a page. This image or word should represent the overall theme of the map and will be important in categorizing subtopics. Once the central point has been established, multiple curved branches will extend out from it and key concepts will be written down along that branch. These topics are called chapter headings or Basic Ordering Ideas (BOIs) when they appear at the first level of branching and should be limited to only a few keywords instead of sentences (Buzan and Griffiths). Every subsequent level past the BOIs will have numerous branches that extend from the previous level. This process keeps going until one reaches the end of the mind map and has a product similar to Figure 1. Some important rules to note when designing a mind map are that the thickness of the branches will decrease for every successive level and that different colors are used to represent each BOI branching (Buzan and Griffiths). Images or doodles should also be incorporated whenever applicable throughout the map to help with association. By

following these steps and requirements, the mind map will optimize information analysis to provide those using the map with numerous cognitive benefits.



**Fig. 1.** Example of mind map from: Brandner, Raphaela. "Why Mind Mapping" *Mind Meister*.

According to Buzan, when creativity is combined with reasoning, the human brain is freed from limitations that standard learning techniques enforce (Buzan 5). This claim is substantiated in a study done by Lawrence E. Murr and James B. Williams, which suggests that utilizing a whole-brain approach while note-taking results in an improvement to intellectual productivity (Murr 418). One specific benefit that mind mapping has is that the amount of information being retained by the user is larger than traditional learning techniques. For instance, in the study of medical students conducted by Farrand, Hussein and Hennessy the memory recall in students increased after two weeks when mind maps were used as a supplemental learning tool (Farrand 430). Mind maps also make it easier for individuals to create comprehensive overviews of a topic, which can lead to improvements in problem solving, negotiations and planning (Buzan

24). For example, a paper was written by a lawyer who stated that mind maps helped allocate assignments to law students in a firm (Stohs 120).

One area in which mind mapping particularly excels is education. Teachers can have their students create mind maps, in which the students identify major and minor themes, either from a lecture or an assigned reading, and draw them out on a piece of paper (Astriani 6). In order for students to get significant improvement in memory and understanding, creative elements like color and images should be incorporated (Tasiwan 125). Mind maps are not only limited to students, however, as teachers can also utilize this tool as a resource when creating lesson plans (Edwards). Another use for mind mapping in education, particularly at the collegiate level, is that it can help students and professors structure an outline for a presentation (Buzan 25).

## 2.2 Natural Language Processing

Language is a form of communication between people through the use of written symbols, sounds or gestures. Languages that are spoken by humans, like English and Spanish, are usually considered to be natural languages because they have been developed through human interaction (Sarkar). Similar to formal languages, like ones that have been designed for computer programming, natural languages are based in syntax and grammatical structure. However, the semantics of any word, phrase or sentence in a natural language has many different factors outside of grammar influencing it. The method in which words obtain their meaning can be split into two categories: lexical semantics and compositional semantics (Sarkar). Lexical semantics involves breaking down a word into morphemes to construct a definition. This method is entirely rule based and is easier for computers to process (Sarkar). Compositional semantics derives the

meaning of a word based on contextual evidence. This could be done in multiple ways such as examining the text surrounding a specific word, analyzing the tone of the overall document or identifying the intent of an author. Compositional meaning is a fundamental part of natural language and is a very difficult piece for computers to evaluate (Sarkar). One of the primary aims of natural language processing is to address this problem.

Natural language processing, or NLP for short, is the intersection between computer science and linguistics. The goal of NLP is to have a computer "comprehend" language in order to extract information that can be used for other purposes (Ghosh). A common example of this is an email service that provides automatic spam identification. These providers are using NLP to analyze an incoming email to determine whether it should be placed in the spam folder (Ghosh). NLP tools that are frequently used can be found in computer programming libraries, like the Natural Language Toolkit (NLTK) or CoreNLP by Stanford University. The method that most NLP technologies use relies on breaking apart text and identifying meaningful sections (Ghosh). For instance, when pre-processing a document, it is common to split up the text into words, phrases, clauses and finally sentences (Sarkar). In order to find every instance of a word in a text, one would use the NLP technique of tokenization ("Natural Language Toolkit."). This process stores every word and punctuation mark into an ordered list and each word in that list can be given a part of speech tag (Ghosh). Once the words have been tagged, parsers, like the RegexpParser from NLTK, can combine these words into different phrases, such as noun and verb phrases (Sarkar). Clauses require two or more phrases and can be split into two classifications, a main clause, which can form a sentence by itself, and a subordinate clause, which requires another clause to form a complete sentence (Sarkar). Another

approach that most NLP researchers take is to clean up the text by removing stop words, like articles, and to normalize text, like changing every instance of "U.S." to "United States" (Ghosh). These methods remove words that provide little understanding and make it easier to apply meaning to a text.

There are many NLP techniques that can be used to clarify the meaning of individual words. One example is called coreference resolution which the Stanford Natural Language Processing Group defines as "the task of finding all expressions that refer to the same entity in a text" ("The Stanford NLP Group."). This clears up ambiguity about what word a pronoun refers to. Another approach is to transform a word into its base form. Two popular methods for this include stemming, which finds the base word by truncating the end of a word when applicable, and lemmatization, which uses a dictionary to find the base word at the cost of program speed (Ghosh). Although understanding the meaning of significant words is valuable, defining the grammar of a sentence or sentences is just as if not more important. One approach to apply meaning to full sentences is called dependency grammar. This method finds the root word, that is, the word that has no dependency on other words in the sentence, and can identify important relationships, such as the nominal subject. Every other word is given a label to describe the relationship it has with either the root word or another word that is modifying the root (Sarkar). An example of the dependency tree produced by the displaCy dependency parser can be seen on the sentence "George Washington became President of the United States in 1789" in Figure 2.

**Fig. 2.** Example of dependency parsed sentence from: displaCy Dependency Visualizer.

Many tasks can be achieved by fundamental natural language processing methods, but as the demand for computers to understand text increases, more advanced technologies, like machine learning (ML), are being used to further the field of computational linguistics. Machine learning can be defined as a set of algorithms that assess and interpret data in order to output predictions or possible solutions to a problem (Theodoridos). Although machine learning is often used in the same context as artificial intelligence, current ML applications are not comprehensive nor advanced enough to be considered full blown AI. Despite this, machine learning has made significant contributions to NLP. For instance, the BERT model, which stands for Bidirectional Encoder Representations from Transformers, performs better than non-ML techniques on various NLP tasks while still being relatively easy to use (Gonzalez-Carvajal 10). Often times, machine learning is being combined with standard NLP methods to produce real word applications. These tasks include question and answering models in which a program answers questions based on text that has been processed, summarization applications where a program condenses a text document into key sentences and machine translation which aims to convert a text from one language to another (Beysolow).

9

**2.3 Related Works**

     Research concerning mind maps and natural language processing was first
analyzed by faculty at Ain Shams University in 2009. The goal of this paper was to use
text analytics methods to generate mind maps when given a text (Abdeen 95). Although
the paper is over a decade old, some of the techniques being used are similar to modern
day NLP techniques. For instance, the researchers were using discourse analysis to clarify
who or what a pronoun was referring to. Today, many NLP libraries include coreference
resolution which achieves the same results. The researchers also used context free
grammars to identify phrases and clauses within a text, resulting in a file that contained
over 600 grammatical rules (Abdeen 97). Just like with coreference resolution, most NLP
libraries offer parsing tools that do not require a programmer to insert context free
grammars. For the time it was published, the results of this study were notable but the
generated mind maps could only evalutate a small amount of the English language
(Abdeen 99).

     There were two more key studies that came out in 2012 in regard to NLP and
mind maps. The first was done by researchers Kudelic, Malekovic and Lovrencic who set
out to create another automatic mind mapping software. Instead of using text files, this
study extracted text from web pages, and many of their techniques would fall under the
category of text analytics, not NLP (Kudelic 124). The next study conducted by
professors in the Computer Science department at Rutgers produced the most
comprehensive and compelling results. The aim of the researchers was to automate the
creation of multilevel mind maps, that is mind maps that contain mind maps within itself
(Elhoseiny 326). The text that was being evaluated was a set of articles about notable

historical figures and the tools that researchers used were NLP techniques. While the subjective accuracy of this study was adequate, the research is over 8 years old which is considered a long time in the field of NLP. There are many techniques that they used which have since been improved upon and there are new techniques that can be used to further enhance the automated mind maps.

Although the research stated above provides a strong foundation for any automated mind mapping system, this thesis aims to use current NLP tools to analyze educational material. Some of the methods that will be used in this project are similar to the related works, but newer tools, like text summarization and improved data visualization, will be used to try and improve the accuracy and design of previous models. It is important to mention that all the studies stated above incorporated images into their systems in order to adhere to standard mind mapping practices. This thesis will not be including images in the final mind map to ensure that more effort is concentrated into text extraction. Another concern with putting images into the mind map is that since the text being analyzed is a college level U.S. history textbook, many of the key concepts that will be placed in the map are hard to define with imagery.

## Chapter 3: Methodology

### 3.1 Preprocessing the Text

Mind maps can be constructed from any source material, but since this research uses natural language processing, history textbooks serve as the ideal educational material to base the mind mapping system on. Most NLP research is based on texts like newspapers, novels, articles and short stories and although history textbooks are not identical to these sources, the primary focus of all these texts is to tell a story to its audience (Bird). Therefore, existing NLP tools will work better with the content inside a history textbook than it would other subjects. Another reason why history is the best subject matter for this research is that it is less likely to contain symbols or sequences that is difficult for NLP tools to manage. For instance, a textbook on calculus is likely to have numerous equations, graphs and mathematical symbols embedded within the text and, although NLP is advancing rapidly, current technologies do not have a method to separate this mathematical notation from the text. The textbook that will be used for this project will be concentrated on U.S. history. Many college students studying in the U.S. are required to take an American history class and would, therefore, find a studying tool based on U.S. history textbooks to be useful with their classes. In addition to this, the researcher of this project is familiar with U.S. history, which makes the results of the mind mapping model easier to evaluate. For example, if the source of this research's mind maps were an economics textbook, it would be harder for the researcher, who has little expertise in economics, to determine if the generated mind map accurately reflects the text.

A digital version of a U.S. history textbook is used in this project since the evaluation of the text is done through computer programming. A PDF file of the textbook is downloaded off the internet and converted into a file format that makes it possible for NLP tools to analyze the text. The format can either be a JavaScript object notation (JSON) file or a standard text file. A JSON file has the same structure as the object data type from JavaScript so when the file gets passed to a web server, the server can interpret the data it contains. For this project, only a text file was used because it makes text extraction easier. After the textbook had been converted, the size of the text being analyzed needed to be established. If the project were to examine an entire textbook or chapter, the resulting mind map would be massive and very difficult to interpret. Since the purpose of mind mapping in this instance is to make it easier for students to comprehend a text, smaller sections of the textbook are used as the mind map source. Therefore, this research only uses subsections of select chapters to generate the mind mapping model to ensure that it is a valuable tool for students.

**3.2 Identifying Key Concepts and Corresponding Properties**

There are two major parts to constructing a mind map, a main idea and the key concepts that branch off the main idea. Given the layout of the textbook being used, the main idea can be easily extracted from the title of the section within a chapter. Similar to the main idea, the first level of key concepts, or the BOIs, can be extracted from the title of each subsection within the chapter section. The subsection title summarizes what the following text consists of in a few words which is ideal for the BOI. To identify the remaining key concepts, the significant portions of the subsection text need to be

identified and isolated. However, before this extraction takes place, the text needs to be cleaned up.

The current state of the text has a large amount of anaphoric references within it. For instance, the two sentences "Still other women accompanied the army as 'camp followers,' serving as cooks, washerwomen, and nurses" and "A few also took part in combat and proved their equality with men through violence against the hated British" refer to a group of women performing actions. If the second sentence were isolated from the first, however, it would not be possible to determine what the words "a few" are in reference to. This can be remedied with an NLP tool called coreference resolution, whose function is to resolve all references to an object. By using this method, we can improve the results of other NLP techniques that extract key concepts from significant portions of the text. Although coreference resolution models have advanced significantly in the past two decades, current models still produce errors and struggle with certain challenges, like cataphora detection (Maslankowska). In this project, a neural network approach, from the NeuralCoref library, is used. This approach was chosen because it produces some of the most accurate results among other models and it is versatile for many different needs (Maslankowska). NeuralCoref also has certain parameters that can be customized by the user to adapt to their specific needs.

Once the coreferences have been resolved, the most important sentences need to be extracted from the text. Instead of spending a considerable amount of time developing a text significance model, this project uses the BERT Extractive Summarizer. This summarizer is defined within a Python library and uses machine learning to implement a pre-trained text summarization model (Liu). It produces high quality summaries in an

optimized amount of time while still being easy to implement (González-Carvajal 8). The NeuralCoref library can be accessed while using the BERT summarizer which makes it ideal for this project. Unfortunately, the text summarizer does not always produce correct results. Since determining accurate summaries is a subjective process, it is harder for algorithms and machine learning techniques to determine exact answers. Therefore, there is room for improvement with the BERT summarizer, but it currently remains as one of the best extractive summarizers (Flannery 26).

After the most important sentences have been extracted, the key concepts for each sentence need to be identified. This is done through dependency parsing which uses the grammatical structure of a sentence to determine dependencies between words. Once the dependency parser has identified the root of the sentence, or the word that conveys the key aspect of the sentence, it labels the relationship other words have to the root word. One of the dependencies the parser recognizes is the nominal subject of a sentence. For this project, every subject that the parser identifies in the sentences of the summary will be labeled as key concept. Then, for every part of the sentence that does not have a direct relationship with the subject, they will be labeled as the properties of the key concept. An example would be the sentence: "At first, Hollywood encountered difficulties in adjusting to the post-World War II environment." Hollywood is identified as the subject of the sentence and every other part of the sentence is a property of Hollywood. Figure 3 shows what the example sentence would look like in mind map format. With the main idea and key concepts for each text established, these items need to be virtually displayed in the style of a mind map.

**Fig. 3.** The sentence -"At first, Hollywood encountered difficulties in adjusting to the post-World War II environment"- in a mind map design made on app.mindmup.com.

## 3.3 Graphically Depicting the Data as a Mind Map

The first step in graphically representing the data is to arrange it into an appropriate format. In this project, the information collected from the previous section is compiled into a JSON file. By using a JSON file, the data is organized into a hierarchical model which is essential to the structure of mind maps. A data visualization tool is needed to extract the text from the JSON file and display the concepts with their respective relationships. The depiction needs to have the main idea located in the center with lines extending from the center to the BOIs. More lines extend out from each BOI until the end of the branch has been met in the JSON file. To accomplish this design, the research uses the D3.js JavaScript library to transform the JSON file into a digital display on a webpage. D3 uses JavaScript to incorporate the data from the processed JSON file with HTML, CSS and Scalable Vector Graphics (SVG), 3 major applications used in web development. The D3 library is ideal for this research because the SVG that displays the

16

mind map has numerous properties that implement design features. For example, the

color and thickness of the branch can be specified, which are two major requirements of

traditional mind maps. The final product is displayed in a webpage that is easily

accessible for students.

## Chapter 4: Conduct of Experiment

### 4.1 Textbook Extraction

This project used the textbook titled *U.S. History* by P. Scott Corbett, Volker Jansen, John M. Lund, Todd Pfannestiel and Paul Vickery. It was originally published in 2014 by OpenStax, a company that specializes in publishing open license textbooks that can be freely downloaded off their webpage as a pdf file (Corbett et al.). As mentioned in Section 3.1, the ideal file format when performing natural language processing is either a text file or a JSON file. Since the textbook was downloaded as a PDF, the Python library called pdfplumber was used to parse the pdf so the text could be transferred to a text file. Pdfplumber works by splitting a pdf into pages and locating every character sequence on the page. One problem pdfplumber has, however, is that it cannot always recognize spaces within a text, especially if the spaces are not always uniform. To fix this issue, the x tolerance, or the amount space between two characters along a horizontal axis, was set to 0.5. This value was accurately able to recognize almost every appropriate space within the textbook. Additionally, pdfplumber only has the capability to recognize text within a pdf, not images. Therefore, every figure, graph and picture within the textbook have been omitted from the text file.

Although the textbook conversion between two files was an automated process, the sections that were chosen for analysis needed to be manually extracted. The main reason for this decision was that chapter sections were chosen based on their subject matter. For testing purposes, these sections needed to cover topics that many who have studied U.S. history would be familiar with. This way, it was easier for the researcher to determine the accuracy of the final product. Size was also taken into account because

some sections had too much information that would cause the mind map to be overly complex. Picking sections based on these criteria was better done manually by the researcher instead of automatically through a computer program. Figure 4 specifies all the sections chosen for analysis. It should be noted that Section 4 is considerably longer than all the other sections. This was done to test how the model would react to a text with a significant amount of detail.

| Section 1: Identity during the American Revolution | Section 2: Indian Removal |
|---|---|
| Section 3: Popular Culture and Mass Media | Section 4: The Rise of Franklin Roosevelt |

**Fig. 4.** Four sections chosen from chapters in *U.S. History*.

The text for each section needed to be filtered in order to remove irrelevant components. For instance, even though the figures were removed when the pdf was converted, the text for each figure was still included throughout the textbook and was subsequently present in the section text files. Additional pieces resided within the section text, like the page numbers, headers, footers and review questions. The removal of these parts was done manually by the researcher since it was more feasible than writing a computer program to address a large number of rules specific only to *U.S. History*. The first paragraph of the subsection text was also removed because it provides a vague

summary of what the section is about to cover but does not provide enough details to be semantically significant to the final product. Furthermore, each section text file was split into subsections that were separated by the corresponding uppercase titles.

**4.2 Subsection Summarization**

For this model, important concepts in each subsection were identified by using the BERT extractive text summarizer. However, in order to optimize the performance of the summarizer, coreference resolution was applied to the text first. This was done through the NeuralCoref package which is a pipeline extension of the open-source NLP python library spaCy ("Neuralcoref"). NeuralCoref uses neural networks and clustering algorithms to resolve coreferences. The main purpose of this package is to locate every mention of an entity and to return data on each mention to the user. It also has another feature where a programmer can resolve the document by replacing every mention in the text with the entity or object it refers to. After utilizing this option, NeuralCoref produced some errors throughout the section text, such as linking a mention to an incorrect entity or only replacing a mention with the first half of the entity's name. To resolve this, the greediness parameter, or the amount of coreference links being made based on a value from 0 to 1, was changed from the default value of 0.5 to 0.4. This number was chosen because it corrected a majority of the errors while still resolving most coreferences. After the coreferences had been specified, the BERT extractive text summarizer was used for each subsection.

The BERT model for extractive text summarization works by embedding every sentence in the text and applying a clustering algorithm to return the sentences that are closest to the cluster centroid (Miller 1). This means that the summarizer is assigning a

numerical value to each sentence and uses these numbers to group similar sentences together. Within each group, a central point is identified and the sentences that are near this point are chosen for the final summarization. Since the summary is extractive and not abstractive, the structure of the sentences that were returned did not differ from how they appeared in the original subsection text. This caused some problems, however, because the summarizer occasionally returned sentences that introduced key concepts but did not give many details surrounding the concepts. For example, the sentence "The HUAC hearings also targeted Hollywood" is identified as being significant because it talks about HUAC, or the House Committee on Un-American Activities, and Hollywood, both of which are important concepts in the section on Popular Culture and Mass Media. This sentence would work in a general summary but for this project's goal, it would be better to choose a sentence that mentions the HUAC's relationship with Hollywood while also providing specific details on it.

To address these issues, the BERT model has some parameters that can be defined in order to specify what sentences should be chosen for the final summary. One of these parameters is called "min_length" and it holds the minimum number of characters a sentence should be in order to be included in the summary. For this project, the value was set to 50 characters while the default is usually 25 characters. This eliminated shorter sentences that included the key concept but had little semantic information regarding the concept. Another parameter that was used is called "ratio" which determines how many sentences should be extracted based on the length of the initial text. While the default of this parameter is set to 0.2, it was increased to 0.3 in this project in order to include more key concepts that were necessary to the understanding of the text. Given these

modifications to the model, a summary of each subsection was created and stored for dependency parsing.

**4.3 Parsing**

The key concepts for each sentence were identified by using the Stanford Core NLP dependency parser. This model identifies each word's part of speech and employs transition-based parsing, or parsing that is gradually developed on a word by word basis (Chen and Manning 741). Once the relationships between each word are identified, they are given a label to describe their function within a sentence. A list of the most common labels is given in Figure 5. Before the subsection summaries were parsed, some stop words, or words that hold no semantic value, were extracted. Therefore, the following words were removed from every sentence: a, an, the, this, that, these, those. This was done because articles and other stop words are typically excluded during notetaking since the notetaker can still interpret meaning without the presence of these words.

root - root
dep - dependent
    aux - auxiliary
        auxpass - passive auxiliary
        cop - copula
    arg - argument
        agent - agent
        comp - complement
            acomp - adjectival complement
            ccomp - clausal complement with internal subject
            xcomp - clausal complement with external subject
            obj - object
                dobj - direct object
                iobj - indirect object
                pobj - object of preposition
        subj - subject
            nsubj - nominal subject
                nsubjpass - passive nominal subject
            csubj - clausal subject
                csubjpass - passive clausal subject
    cc - coordination
    conj - conjunct
    expl - expletive (expletive "there")
    mod - modifier
        amod - adjectival modifier
        appos - appositional modifier

advcl - adverbial clause modifier
det - determiner
predet - predeterminer
preconj - preconjunct
vmod - reduced, non-finite verbal modifier
mwe - multi-word expression modifier
    mark - marker (word introducing an advcl or ccomp
advmod - adverbial modifier
    neg - negation modifier
rcmod - relative clause modifier
quantmod - quantifier modifier
nn - noun compound modifier
npadvmod - noun phrase adverbial modifier
    tmod - temporal modifier
num - numeric modifier
number - element of compound number
prep - prepositional modifier
poss - possession modifier
possessive - possessive modifier ('s)
prt - phrasal verb particle
parataxis - parataxis
goeswith - goes with
punct - punctuation
ref - referent
sdep - semantic dependent
    xsubj - controlling subject

**Fig. 5.** List of grammatical relations specified by the Stanford Dependency Parser from: Marneffe and

Manning "Stanford typed dependencies manual."

Once the modified sentences were passed through the dependency parser, every

word that was given a subject label, like the nominal subject (nsubj) or clausal subject

(csubj), was identified within each sentence. Every word that had a direct relationship to

the subject was located and added to a key concept variable. The rest of the sentence was

added to a variable that held the attributes of a key concept. This typically ended up

splitting the sentence into one noun phrase which held the subject and one verb phrase.

Based on where they resided within the section text, both the key concept and attribute

variables were stored in a hierarchical format. In some instances, multiple sentences

contained the same key concept. To reduce unnecessary repetition, the attribute variables

23

were combined and coupled with only one key concept variable. Ultimately, the hierarchy that was transferred to a JSON file started with the main concept, which was the chapter section title, then the basic ordering idea (BOI), which was the subsection header, and ended with the subsection key concepts and their attributes. An example of the format is given in Figure 6.



**Fig. 6.** Example of the JSON file hierarchy.

## 4.4 Data Visualization

D3.js is a JavaScript library that utilizes the Document Object Model of HTML files to manipulate structured data (Bostock, "Data-Driven Documents"). HTML, which stands for Hypertext Markup Language, is the language used to display documents in a web browser. It is often linked to CSS, or Cascading Style Sheets, and JavaScript files to give the resulting webpage more functionality. Therefore, by using specified data with the D3 library, an HTML file will display the data based on the JavaScript written inside it.

**Fig. 7.** Example of the radial tidy tree from: Bostock "Radial Tidy Tree."

The D3 website has many examples of how the library can be used to visualize various datasets. For this project, the radial tidy tree, as seen in Figure 7, was used because it contains a central point and various nodes extending from that point which is similar to the structure of a mind map. This example requires hierarchical data and uses the Reingold-Tilford "tidy" algorithm to represent the diagram as an SVG (Bostock, "Radial Tidy Tree"). By passing the JSON file with the processed text to the D3 tidy tree

implementation, the data was displayed through a web browser. Although the output

matched the radial tidy tree example, more adjustments needed to be made in order to

display the diagram as a mind map. Specifically, the text at the end of each node was

only being displayed on a single line instead of wrapping to a new line. This cut off parts

of the text that extended past the width of the SVG and caused some node texts to overlap

with each other. To fix this, the text appended to each node was changed from a text

element to a Foreign Object element. The text was then given a maximum width that it

could not exceed which resulted in the text wrapping around itself.

**Chapter 5: Results**

**5.1 Overall Results and Analysis**

The final result of this mind mapping system was four mind map inspired graphs that were generated based on their respective text section. The graph for Section 1 can be seen in Figure 8 and provides the best visual representation of all the sections analyzed. This was mainly due to the organization of Section 1's text because it evenly split its content into four subsections.

Because of the BERT Text Summarizer, each graph does a pretty accurate job at representing the major points of its corresponding text. For example, in the graph of Section 3, the branch titled "Rocking Around the Clock" does a great job expanding on almost every important concept mentioned in the first subsection of Section 3, which can be found in the appendix. There were some minor and major concepts that the summarizer missed, however, because the model had a set maximum for the number of sentences it could choose. Additionally, the BERT Text Summarizer was further limited because it is an extractive system which means it did not rephrase the sentences in the text, it only pulled the ones it considered significant.

**6.4 Identity during the American Revolution**

- **WOMEN**
  - women
    - In colonial America, shouldered enormous domestic and child-rearing responsibilities.
    - were also expected to provide food and lodging for armies and to nurse wounded soldiers.
    - formed Ladies Association of Philadelphia and led fundraising drive to provide sorely needed supplies to Continental Army.
    - Still accompanied army as "camp followers," serving as cooks, washerwomen, and nurses.
    - also took part in combat and proved their equality with men through violence against hated British.
  - Esther DeBerdt Reed of Philadelphia, wife of Governor Joseph Reed,
  - other women
  - A few

- **PATRIOTS**
  - inflation
    - By 1781, was such 146 Continental dollars were worth only one dollar in gold.
  - British
    - Whereas could pay in gold and silver, American forces relied on paper money, backed by loans obtained in Europe.
  - shortage of supplies
    - After initial burst of enthusiasm in 1775 and 1776, became acute in 1777 through 1779, as Washington's difficult winter at Valley Forge demonstrates.
  - The American revolutionaries (also called Patriots or Whigs)
    - came from many different backgrounds and included merchants, shoemakers, farmers, and sailors.

- **LOYALISTS**
  - Historians
    - disagree on what percentage of colonists were Loyalists; estimates range from 20 percent to over 30 percent.
  - themselves
    - could and preferring not to engage in struggle.
  - revolutionaries
    - During war, imprisoned William in Connecticut; however, he remained steadfast in he allegiance to Great Britain and moved to England after Revolution.
  - Another sizable group of Loyalists
    - went to British West Indies, taking Another sizable group of Loyalists slaves with Another sizable group of Loyalists.

- **SLAVES AND INDIANS**
  - white revolutionaries
    - In new United States, largely reinforced racial identity based on skin color.
  - Revolution
    - more than two refused to let slaves serve in army, although he did allow free blacks.
    - Washington, owner of hundred slaves during Revolution,
  - some slaves who fought for Patriot cause
    - While received some slaves who fought for Patriot cause freedom, revolutionary leaders — unlike British — did not grant such slaves freedom as matter of course.
  - Between ten and twenty thousand slaves gained their freedom because of Revolution; arguably created largest slave uprising and greatest emancipation until Civil War.
  - Indeed, despite their class and ethnic differences, stood mostly united in white revolutionaries hostility to both blacks and Indians.

**Fig. 8.** Mind map inspired graph for Section 1.

The dependency parser performed pretty well in all of the graphs and accurately split the key concepts from their respective attribute in most cases. This is demonstrated best in Figure 9 where the subjects of each sentence are identified as a key concept and the verb phrases linked to the subjects are correctly represented. There were a few mistakes, however, where some words were excluded from the key concept and were either incorrectly added to an attribute or left off entirely. An example of this can be found in Figure 8 under the "Slaves and Indians" branch where the second key concept "Washington, owner of hundred slaves during revolution" is missing the phrase "more than two" which changes the context of the concept. Mistakes like this were caused by either an error from the parser, an error in the code that split the subject from the rest of the sentence, or both.



**Fig. 9.** Example of the separation between key concepts and their attributes from Section 2.

One recurring problem in most of the graphs was that some key concepts referred to vague entities. For instance, in Figure 10 the key concept "Many" does not give the

reader an idea of who exactly regarded the music as a threat. Another case can be found in Figure 8 where the concept "themselves" under the "Loyalists" branch was not only given a vague label, but it also referenced an incorrect entity. These inaccuracies were caused by coreference resolution which would sometimes fail to recognize a reference as a mention of an entity, or it would match a mention to an incorrect entity. The only graph in which this was not a problem was the one that represented Section 2. One explanation for why this was the case is that the number of words in each section had an impact on the performance of coreference resolution. Section 2 was the second shortest in length and section 3, the shortest one, contained only one coreference error. Sections 1 and 4, the two longest sections, each contained two coreference mistakes which implies that coreference resolution performs worse on longer texts.



**Fig. 10.** Error with coreference resolution from Section 3.

Data visualization worked well with some of the sections but lacked the capabilities to represent others. The radial design feature of D3.js proved to be very useful and created a visually appealing structure that imitated a mind map. Furthermore, the implementation of colors to separate the branches improved the readability of the graph. The best representation of both of these features can be seen in Figure 8. One

30

element that D3.js struggled with was the text placement. The width and placement of all the texts had to be manually specified because there was no responsive way to structure them. This was particularly problematic when large portions of the section text needed to be placed on the graph, as seen in Figure 11. Additionally, the thickness of each branch could not be decreased for every subsequent level, which is a common requirement for mind maps.



**Fig. 11.** Overlapping texts from Section 4.

## 5.2 Results of Testing a Large Text

As mentioned in the Conduct of Experiment, this size of Section 4 was considerably bigger than all the other sections. This was done to test how this system would work on a larger text. The resulting graph, as seen in Figure 12, shows that the Section 4 performed significantly worse than the other sections. In terms of content, there was a substantial amount of material that was displayed, with one subsection containing 7 key concepts. Moreover, some of these key concepts should not have been considered significant enough to add to the graph, like the phrase "how to take action on beliefs" under the "The Interregnum" branch in Figure 12. Regarding visualization, the design

31

implementations of the graph caused Section 4 to be unreadable in certain areas. For example, the attributes of the concept "Franklin Roosevelt" in the upper right corner overlap with each other and make it impossible to understand what each phrase is saying. This could be corrected by using a different data visualization library, but even if all the concepts were separated and legible, the amount of information on one page would be too overwhelming for a reader to interpret which defeats the purpose of creating a mind map inspired graph. Ultimately, this system works best with smaller sections.

THE INTERREGNUM

INAUGURATION DAY: A NEW BEGINNING

THE ELECTION OF FRANKLIN ROOSEVELT

26.1 The Rise of Franklin Roosevelt

Roosevelt

who

one observer in crowd

Upon hearing Roosevelt inaugural address, later commented, " Any man who can talk like in times like is worth every ounce of support true American has "

Borrowing wartime analogy provided by Moley, served as Roosevelt speechwriter at time, Roosevelt called upon all Americans to assemble and fight essential battle against forces of economic depression.

Bathed in sunlight, delivered one of most famous and oft-quoted inaugural addresses in history.

rode in open car along with outgoing president Hoover, facing public, as Roosevelt made Roosevelt way to U.S Capitol.
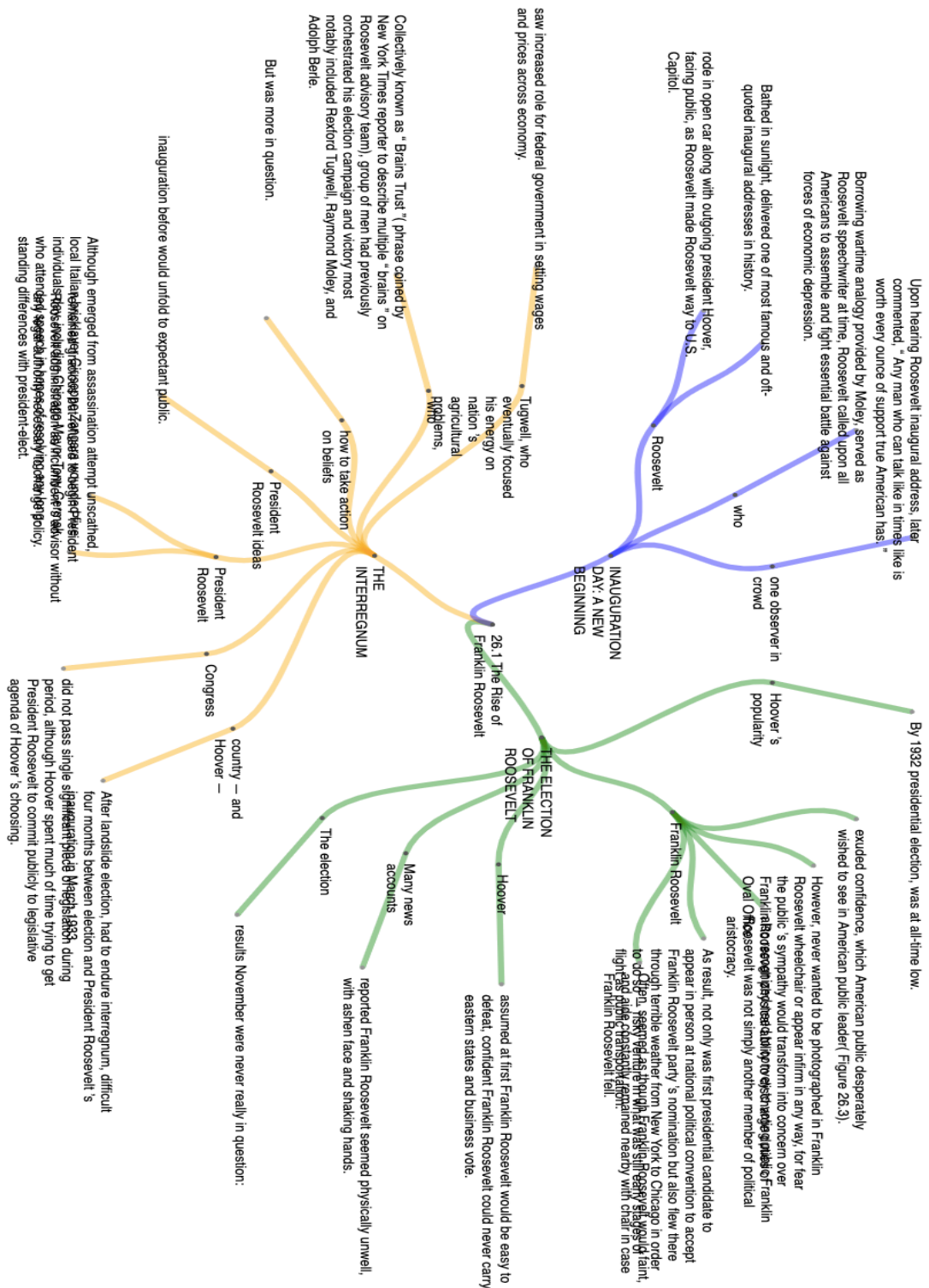
saw increased role for federal government in setting wages and prices across economy.

Tugwell, who eventually focused his energy on nation 's agricultural problems, RFC

how to take action on beliefs

President Roosevelt ideas

President Roosevelt

President Congress

country — and Hoover —

Collectively known as " Brains Trust " ( phrase coined by New York Times reporter to describe multiple " brains " on Roosevelt advisory team), group of men had previously orchestrated his election campaign and victory most notably included Rexford Tugwell, Raymond Moley, and Adolph Berle.

But was more in question.

inauguration before would unfold to expectant public.

Although emerged from assassination attempt unscathed,

Hoover 's popularity

By 1932 presidential election, was at all-time low.

Franklin Roosevelt

Hoover

Many news accounts

The election

results November were never really in question:

exuded confidence, which American public desperately wished to see in American public leader( Figure 26.3).

However, never wanted to be photographed in Franklin Roosevelt wheelchair or appear infirm in any way, for fear the public 's sympathy would transform into concern over Franklin Roosevelt was not simply another member of political aristocracy.

As result, not only was first presidential candidate to appear in person at national political convention to accept Franklin Roosevelt party 's nomination but also flew there through terrible weather from New York to Chicago in order to

assumed at first Franklin Roosevelt would be easy to defeat, confident Franklin Roosevelt could never carry eastern states and business vote.

reported Franklin Roosevelt seemed physically unwell, with ashen face and shaking hands.

After landslide election, had to endure interregnum, difficult four months between election and President Roosevelt 's inauguration in March 1933, during period, although Hoover spent much of time trying to get President Roosevelt to commit publicly to legislative agenda of Hoover 's choosing.

did not pass single significant piece of legislation during

**Fig. 12.** Mind map inspired graph for Section 4.

## Chapter 6: Conclusion and Future Work

This research aimed to generate mind map inspired graphs based on sections from a U.S. history textbook in order to decrease the amount of workload required from users. This was done using natural language processing techniques, such as extractive text summarization and dependency parsing, along with the D3.js data visualization library. Overall, the result of this research was successful in creating a mind mapping system that requires minimal input from targeted users, like students and teachers. Some of the final graphs contained errors that were largely due to the limitations of current natural language processing technologies. Additionally, some of the sections contained too much text, which could not be clearly represented in the graphs using D3.js. These shortcomings can be resolved in the future by using more advanced NLP and visualization tools.

The most promising prospect of this research is that, with further development, it can be applied to fields outside U.S. history. This subject was chosen because it worked best with what the researcher wanted to produce, but a new system that either analyzes a specified subject outside of history or an array of subjects could be created using some parts of this project. However, this system was designed with one textbook in mind, which means that some features will only work if the text being analyzed has the same structure as the textbook. In order to work with different texts, future researchers would need to address these challenges.

One key point from this research is that the graphs produced are only inspired by mind maps. They do not satisfy all the requirements set forth by Tony Buzan, the creator of mind maps, and, therefore, cannot be considered as such. Future work can expand on

what has already been accomplished by including features that this model currently lacks. For example, most mind maps contain either pictures or doodles that help the user associate a key concept with an image. This project did not include images in the final graphs because it would not have been feasible to complete everything within the time given. Therefore, future researchers can explore different methods for adding pictures to this system's current output. This could be done by extracting images off the internet when searching for a key word or by using the figures that were presented within the textbook.

The long-term prospect of this project is that the accuracy can improve with the advancement of natural language processing techniques. For instance, an abstractive text summarizer might be a better alternative to an extractive version because it is not limited by the structure of sentences within the target text. An abstractive model rephrases the concepts it considers to be significant into new sentences and organizes them to form one coherent summary, which could improve the accuracy of the graphs in a project like this. Unfortunately, research into abstractive models has not yet reached the point where it can be reasonably implemented into other projects. When this technology is readily available, however, it can be used to further the results of this research.

Works Cited

Abdeen, Mohammad, et al. "Direct automatic generation of mind maps from text with M
2 Gen." 2009 IEEE Toronto International Conference Science and Technology for
Humanity (TIC-STH). IEEE, 2009.

Astriani, Dyah, et al. "Mind mapping in learning models: A tool to improve student
metacognitive skills." International Journal of Emerging Technologies in
Learning (iJET) 15.6 (2020): 4-17.

Beysolow II, Taweh. Applied Natural Language Processing with Python: Implementing
Machine Learning and Deep Learning Algorithms for Natural Language
Processing. Apress, 2018.

Bird, Steven, et al. "2. Accessing Text Corpora and Lexical Resources." *Natural
Language Processing with Python*. 3rd., 2019.

Bostock, Mike. "D3.Js - Data-Driven Documents." D3, 2012, http://d3js.org/.

Bostock. "Radial Tidy Tree." D3, 27 Aug. 2020, observablehq.com/@d3/radial-tidy-tree.

Buzan, Tony, and Chris Griffiths. Mind Maps for Business 2nd edn: Using the ultimate
thinking tool to revolutionise how you work. Pearson UK, 2013.

Buzan, Tony. *Mind Map Handbook: The Ultimate Thinking Tool*. E-book, Harper
Collins, 2004.

Chen, Danqi, and Christopher D. Manning. "A fast and accurate dependency parser using
neural networks." Proceedings of the 2014 conference on empirical methods in
natural language processing (EMNLP). 2014.

Corbett, P. Scott, et al. U.S. History by OpenStax. 1st ed., XanEdu Publishing Inc, 2014,
https://openstax.org/details/books/us-history.

Edwards, Sarah, and Nick Cooper. "Mind mapping as a teaching resource." The clinical

    teacher 7.4 (2010): 236-239.

Elhoseiny, Mohamed, and Ahmed Elgammal. "English2mindmap: An automated system

    for mindmap generation from english text." 2012 IEEE International Symposium

    on Multimedia. IEEE, 2012.

Farrand, Paul, Fearzana Hussain, and Enid Hennessy. "The efficacy of the 'mind

    map' study technique." Medical education 36.5 (2002): 426-431.

Flannery, Jeremiah. "Using NLP to Generate MARC Summary Fields for Notre Dame's

    Catholic Pamphlets." International Journal of Librarianship 5.1 (2020): 20-35.

Ghosh, Sohom, and Dwight Gunning. *Natural Language Processing Fundamentals*. 1st

    edition., *Packt Publishing*, 2019.

González-Carvajal, Santiago, and Eduardo C. Garrido-Merchán. "Comparing BERT

    against traditional machine learning text classification." arXiv preprint

    arXiv:2005.13012 (2020).

Huggingface. "Neuralcoref." GitHub, 8 Apr. 2019, github.com/huggingface/neuralcoref.

Kudelić, Robert, Mirko Maleković, and Alen Lovrenčić. "Mind map generator software."

    2012 IEEE International Conference on Computer Science and Automation

    Engineering (CSAE). Vol. 3. IEEE, 2012.

Maslankowska, Marta, and Pawel Mielniczuk. "Most Popular Coreference Resolution

    Frameworks." *Medium*, Towards Data Science, 22 Jan. 2021,

    towardsdatascience.com/most-popular-coreference-resolution-frameworks-

    574ba8a8cc2d

Mento, Anthony J., Patrick Martinelli, and Raymond M. Jones. "Mind mapping in
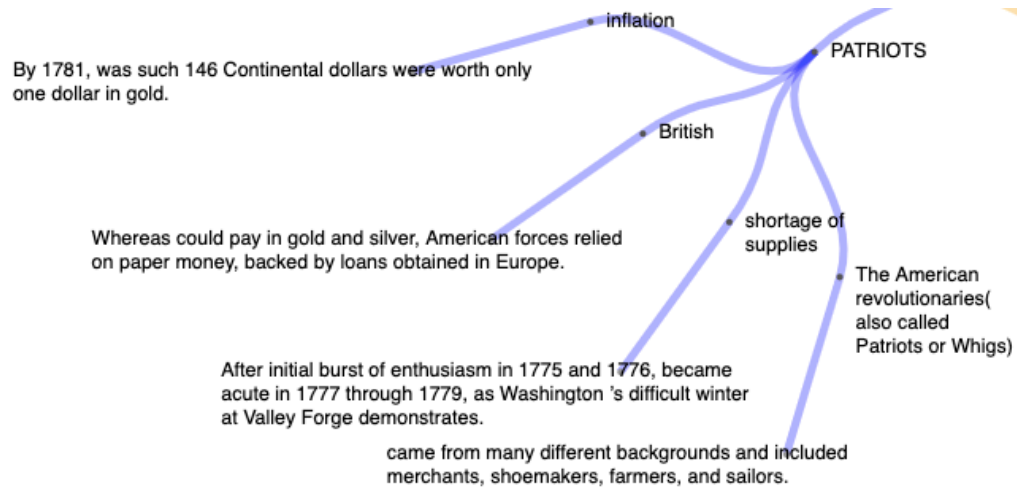
executive education: Applications and outcomes." Journal of Management

Development (1999).

Miller, Derek. "Leveraging BERT for extractive text summarization on lectures." arXiv

preprint arXiv:1906.04165 (2019).

"Mind Mapping or Concept Mapping" *Florida Atlantic University Libraries*,

https://libguides.fau.edu/mindmapping

Murr, Lawrence E., and James B. Williams. "Half-brained ideas about education:

thinking and learning with both the left and right brain in a visual culture."

Leonardo (1988): 413-419.

"Natural Language Toolkit." NLTK Project, Apr. 2020, https://www.nltk.org/

Spencer, Julie R., Kelley M. Anderson, and Kathryn K. Ellis. "Radiant thinking and the

use of the mind map in nurse practitioner education." Journal of Nursing

Education 52.5 (2013): 291-293.

Stohs, Brett C. "Oh What a Tangled Web We Weave: Mind Mapping as Creative Spark

to Optimize Transactional Clinic Assignments." NYL Sch. L. Rev. 61 (2016):

119.

Tasiwan, Tasiwan. "Transformation of the Students' Inquiry Capability Through

Mindmap Educative by Using Game Observation Normatively (Megono)

Learning Model." Jurnal Pendidikan IPA Indonesia 5.1: 123-133.

"The Stanford NLP Group." The Stanford Natural Language Processing Group,

https://nlp.stanford.edu/projects/coref.shtml

Theodoridis, Sergios. Machine learning: a Bayesian and optimization perspective.

Academic press, 2015.

**Excerpts**

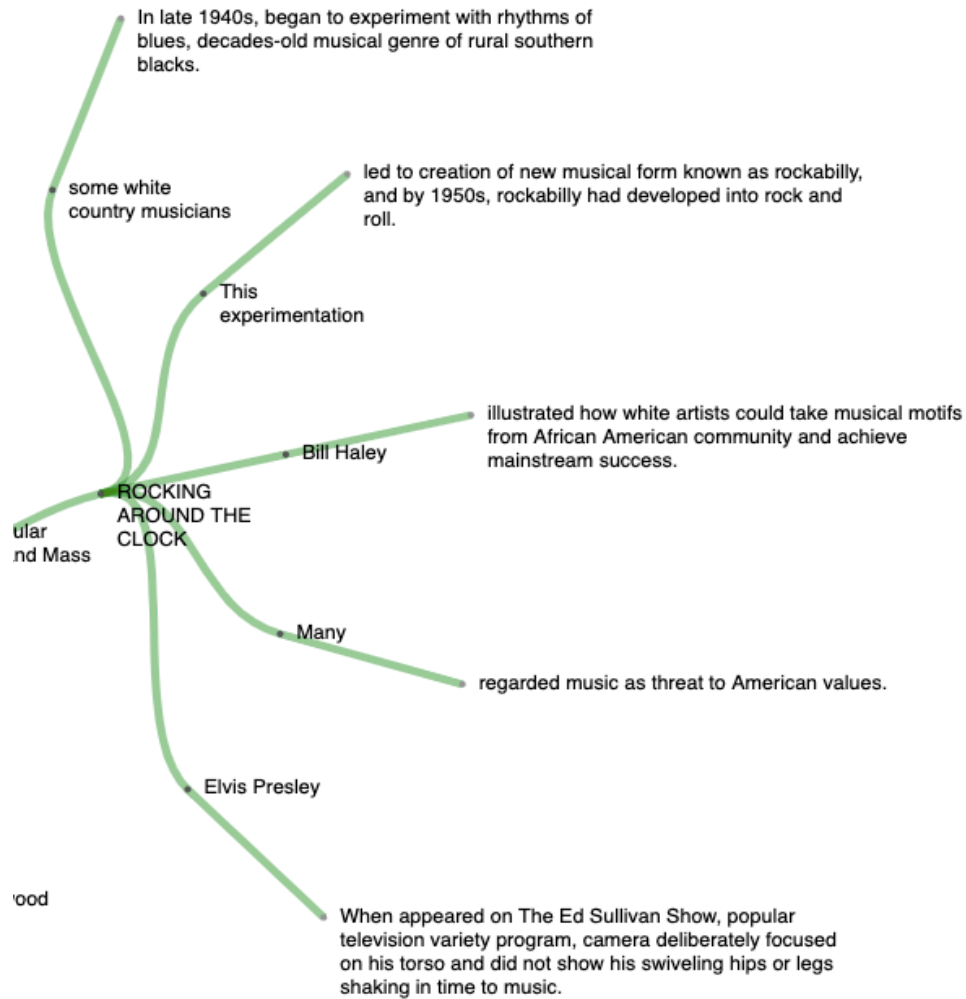Text: Section 1
Subsection: PATRIOTS



The American revolutionaries (also called Patriots or Whigs) came from many different backgrounds and included merchants, shoemakers, farmers, and sailors. What is extraordinary is the way in which the struggle for independence brought a vast cross-section of society together, animated by a common cause. During the war, the revolutionaries faced great difficulties, including massive supply problems; clothing, ammunition, tents, and equipment were all hard to come by. After an initial burst of enthusiasm in 1775 and 1776, the shortage of supplies became acute in 1777 through 1779, as Washington's difficult winter at Valley Forge demonstrates.

Funding the war effort also proved very difficult. Whereas the British could pay in gold and silver, the American forces relied on paper money, backed by loans obtained in Europe. This first American money was called Continental currency; unfortunately, it quickly fell in value. "Not worth a Continental" soon became a shorthand term for something of no value. The new revolutionary government printed a great amount of this paper money, resulting in runaway inflation. By 1781, inflation was such that 146 Continental dollars were worth only one dollar in gold. The problem grew worse as each former colony, now a revolutionary state, printed its own currency.

(From https://openstax.org/details/books/us-history)

Text: Section 3
Subsection: ROCKING AROUND THE CLOCK

In late 1940s, began to experiment with rhythms of blues, decades-old musical genre of rural southern blacks.

some white country musicians

led to creation of new musical form known as rockabilly, and by 1950s, rockabilly had developed into rock and roll.

This experimentation

illustrated how white artists could take musical motifs from African American community and achieve mainstream success.

Bill Haley

ROCKING AROUND THE CLOCK

ular nd Mass

Many

regarded music as threat to American values.

Elvis Presley

ood

When appeared on The Ed Sullivan Show, popular television variety program, camera deliberately focused on his torso and did not show his swiveling hips or legs shaking in time to music.

In the late 1940s, some white country musicians began to experiment with the rhythms of the blues, a decades-old musical genre of rural southern blacks. This experimentation led to the creation of a new musical form known as rockabilly, and by the 1950s, rockabilly had developed into rock and roll. Rock and roll music celebrated themes such as young love and freedom from the oppression of middle-class society. It quickly grew in favor among American teens, thanks largely to the efforts of disc jockey Alan Freed, who named and popularized the music by playing it on the radio in Cleveland, where he also organized the first rock and roll concert, and later in New York.

The theme of rebellion against authority, present in many rock and roll songs, appealed to teens. In 1954, Bill Haley and His Comets provided youth with an anthem for their rebellion- "Rock Around the Clock" (Figure 28.14). The song, used in the 1955 movie Blackboard Jungle about a white teacher at a troubled inner-city high school, seemed to be calling for teens to declare their independence from adult control.

Haley illustrated how white artists could take musical motifs from the African American community and achieve mainstream success. Teen heartthrob Elvis Presley rose to stardom doing the same. Thus, besides encouraging a feeling of youthful rebellion, rock and roll also began to tear down color barriers, as white youths sought out African American musicians such as Chuck Berry and Little Richard (Figure 28.14). While youth had found an outlet for their feelings and concerns, parents were much less enthused about rock and roll and the values it seemed to promote. Many regarded the music as a threat to American values. When Elvis Presley appeared on The Ed Sullivan Show, a popular television variety program, the camera deliberately focused on his torso and did not show his swiveling hips or legs shaking in time to the music. Despite adults' dislike of the genre, or perhaps because of it, more than 68 percent of the music played on the radio in 1956 was rock and roll.

(From https://openstax.org/details/books/us-history)