

## What questions can humanists and social scientists find answers to with the computational tools in the NLP Suite?

### Contents

Preparing your documents for NLP algorithms .....	2
How many documents do you need? .....	3
What kind of documents do you need? .....	3
What's in your documents? A broad view .....	3
Corpus statistics .....	3
What does your corpus talk about? What are the topics? .....	3
What's in your documents? A closer look .....	3
Syntactic relations between words .....	3
Grammatical analysis of text .....	4
Subject-Verb-Object (SVO) extraction and visualizations (the 5Ws of narrative) .....	4
Pronominal anaphora and coreference .....	4
Semantic relations between words .....	4
Searching documents .....	5
Search words in documents .....	5
Search the CoNLL table .....	5
Search documents for n-grams and co-occurrences .....	5
The NLP Suite has an n-grams/co-occurrences viewer, s .....	5
Search with knowledge-graph tools: DBpedia, YAGO, WordNet .....	5
Getting information from the web about organizations, people, locations .....	6
Sentiments and emotions .....	6
Writing style .....	7
N-grams, Co-occurrences, N-grams/co-occurrences viewer .....	7
Verb modality, tense, and voice .....	8
Nominalization .....	8
Sentence complexity .....	8
Text readability .....	8
Vocabulary .....	8
Unusual words/misspelled words .....	8

Short/long words, vowel words .....	8
Vocabulary richness.....	8
Abstract/concrete vocabulary.....	8
Objective/subjective writing .....	8
In which language are texts written in? .....	8
Who wrote the text? Man or woman?.....	8
Is there dialogue in the text? .....	8
Visualizing words .....	9
Wordclouds .....	9
Sankey plot.....	9
Network graphs.....	10
HTML files .....	10
Pin maps heat maps.....	10
Sunburst pie charts .....	10
Treemap .....	11
Animated time plot .....	11
Box plots .....	11
Bar charts, pie charts, line charts .....	11

<b>Humanists' and social scientists' questions</b>	<b>Computer scientists' answers and available NLP Suite Tools</b>
<b>Preparing your documents for NLP algorithms</b>	

<p><i>How many documents do you need?</i></p> <p>Most NLP tools will allow you to work with one or several documents in a directory (your <i>corpus</i>). Some NLP tools, however, will only work with an input corpus of many documents (e.g., topic modeling or shape of stories).</p> <p><i>What kind of documents do you need?</i></p> <p>NLP tools typically take .txt files in input. There are, however, several algorithms in the NLP Suite that will allow you to convert files from one type to the other (e.g., docx or rtf to txt).</p>	<p><i>Under Pre-Processing Tools</i> select the tool that you need among the wide range of available tools.</p>
<p><b>What's in your documents? A broad view</b></p>	
<p>Is there a way to get basic <b>statistics about your corpus</b>, i.e., the set of documents you are studying? For example, the number of documents, the number of words per document, the sentence length per document?</p> <p><i>Corpus statistics</i></p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Corpus/document(s) statistics (Sentences, words, lines)</li> </ul>
<p><i>What does your corpus talk about? What are the topics?</i></p>	<p><i>Under CORPUS analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Topic modeling (via MALLET &amp; Gensim)</li> </ul>
<p><b>What's in your documents? A closer look</b></p>	
<p><i>Syntactic relations between words</i></p>	
<p>Can computational tools tell you <b>which words in a text are nouns, verbs, or adjectives?</b></p> <p>Yes. There are basic tools called <b>parsers</b> that do precisely that, and more and with a high degree of accuracy (over 90%). They can even give the syntactical relation between words, whether verbs are in the infinitive form, gerundive, past, present, passive or active, whether nouns are singular or plural subjects or objects, and, again, active or passive.</p>	<p>The NLP Suite has several different options for parsers: spaCy, Stanford CoreNLP, Stanza, each with slightly different characteristics of accuracy or speed.</p> <p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p>

	<ul style="list-style-type: none"> <li>Parsers &amp; annotators (BERT, CoreNLP, spaCy, Stanza)</li> </ul>
<p>Can they even tell you <b>which adjectives are used for which nouns</b>? Like in a study of folktales, which adjectives are used for princesses and princes?</p> <p><i>Grammatical analysis of text</i></p> <p>Definitely. Parsers produce in output a CoNLL table that contains all the information required to answer those questions.</p>	<p>In the NLP Suite, the CoNLL table analyzer algorithms will allow you to explore these questions.</p> <p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>CoNLL table analyzer - Search the CoNLL table</li> </ul>
<p>Who does what, where, when, and why? Can NLP help to identify the <b>5Ws of stories</b>?</p> <p><i>Subject-Verb-Object (SVO) extraction and visualizations (the 5Ws of narrative)</i></p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>SVO (Subject-Verb-Object) extractor &amp; visualization</li> </ul>
<p>There is a problem in linguistics known as <b>pronominal anaphora</b> where in such a sentence as “John hurried home. He said he was ill,” where both “he” pronouns (the anaphoric expression) refer to “John” (the antecedent). Can computational tools link pronouns to their nouns (anaphoric expressions to their antecedents)?</p> <p><i>Pronominal anaphora and coreference</i></p> <p>In NLP, this problem is known as <b>coreference resolution</b>. Unfortunately, the accuracy of these algorithms is not as high as for parsers, even for neural network approaches (around 70%?).</p>	<p>There are two points of entry to coreference resolution the NLP Suite and with several different options available.</p> <p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>CoreNLP annotator - coreference (pronominal)</li> <li>Parsers &amp; annotators (BERT, CoreNLP, spaCy, Stanza)</li> </ul>
<p><i>Semantic relations between words</i></p>	
<p><b>Which words come together in texts?</b> Are women/girls mostly associated with adjectives or nouns of physical beauty, family relations, and low-paying occupations in the care industry (waiters, teachers, nurses)? Is the opposite true for men and boys?</p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p>

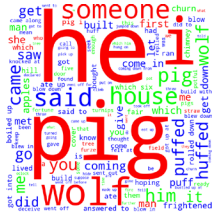
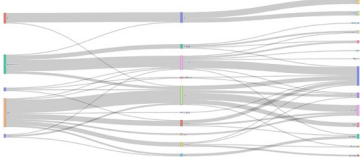
	<ul style="list-style-type: none"> <li>Word embeddings (Word2Vec) (via BERT &amp; Gensim)</li> </ul>
Do <b>nouns and verbs cluster</b> together in specific groups (e.g., verbs of movement or verbs of communication)?	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>WordNet</li> </ul>
<i>Searching documents</i>	
<p>Besides these powerful NLP tools of syntactic and semantic analysis, are there simpler search tools that would allow you to zoom into a text for specific words or combinations of words?</p> <p><b>Search words in documents</b></p> <p>The NLP Suite offers a number of options for searching documents with specific words in mind. Each of these options will allow you to extract the sentence (or a contour of sentences) that contain the search words. And if you work with a large corpus, you can even create a subsample of files that contain the search words for separate analyses.</p> <p><b>Search the CoNLL table</b></p> <p>You can search the CoNLL table, finding, perhaps, all the adjectives used to characterize a specific noun (e.g., prince or princess).</p> <p><b>Search documents for n-grams and co-occurrences</b></p> <p>The NLP Suite has an n-grams/co-occurrences viewer, similar to Google Ngram Viewer (<a href="https://books.google.com/ngrams/">https://books.google.com/ngrams/</a>) but applied to your own corpus, rather than Google books. The viewer can search for words or combinations of words (e.g., “love letter”) and plot them overtime if filenames embed a date (e.g., The New York Times_12-01-1992).</p> <p><b>Search with knowledge-graph tools: DBpedia, YAGO, WordNet</b></p> <p>You can search your corpus for words found in what are known as knowledge-graph tools: DBpedia or YAGO or WordNet, the Princeton University semantic database. Do the</p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>Search (ALL options GUI)</li> </ul>

documents contain words in DBpedia or YAGO ontology classes (e.g., emotions, writer) or in WordNet word categories (“synsets”, e.g., verbs of movement, such as go, come, run)?	
<b>Getting information from the web about organizations, people, locations</b>	
If documents in your corpus mention proper names (people, organizations, locations) can you find <b>information on the web</b> automatically? For instance, when and where someone was born or their occupation?	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Knowledge graphs: DBpedia &amp; YAGO</li> </ul>
<b>Sentiments and emotions</b>	
Which <b>sentiments and emotions</b> do writers express in their writing? Positive, negative? How do they express these emotions rhetorically and linguistically?	<p>The NLP Suite provides several options to address these questions. To find all the options</p> <p><i>Under CORPUS analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Sentiments/emotions (ALL options GUI)</li> </ul> <p>The GUI will give you access to all the various options available in the NLP Suite to find answers to those questions. You can also select</p> <ul style="list-style-type: none"> <li>• Style analysis (ALL options GUI)</li> </ul> <p>Then tick the Vocabulary analysis option and select:</p> <p>Punctuation as figures of pathos (!?)</p>

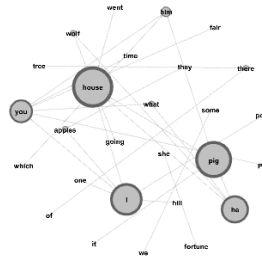
	<p>Or the Repetition tools (repetition is also a figure of pathos).</p> <p>Sentiment analysis (neural network/tensor options: BERT, spaCy, Stanford CoreNLP, Stanza)</p> <p>Sentiment analysis (dictionary options: ANEW, hedonometer, SentiWordNet, VADER)</p> <p>You can also use YAGO's emotion ontology class or WordNet (NOUN feeling VERB emotion) to search for emotion words used in the texts</p>
And when a single story is viewed in the context of hundreds or thousands of other stories, do these stories cluster together in specific <b>story shapes</b> , as Kurt Vonnegut would have it? A man in a hole, from rags to riches or riches to rags? Happy ending or a bummer? Comedy or tragedy?	<p><i>Under CORPUS analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Shape of stories</li> </ul>
<b>Writing style</b>	
Which <b>style</b> do different authors use in their writing? Can the NLP Suite help highlight general characteristics of an author's style?	<p>The NLP Suite provides a wide range of options to find answers to these questions.</p> <p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Style analysis (ALL options GUI)</li> </ul>
Are there words and combinations of words that appear more frequently in texts? How do they differ across different authors?	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• N-grams (word &amp; character)</li> </ul>
<b>N-grams, Co-occurrences, N-grams/co-occurrences viewer</b>	

<p>Computationally, these objects are known as n-grams: unigram for single words, bigrams for combinations of two words, trigrams for three words, etc. (although it generally makes little sense to go beyond trigrams).</p>	<ul style="list-style-type: none"> <li>• N-grams/co-occurrences viewer</li> </ul>
<p>Is there a consistent use of passive and active verb forms, of present and past verb tenses, of some subjects always portrayed in active or passive forms, acknowledging or denying agency?</p> <p><i>Verb modality, tense, and voice</i> <i>Nominalization</i></p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• CoNLL table analyzer - Search the CoNLL table</li> <li>• Nominalization</li> </ul>
<p>How <b>readable</b> is a text? How <b>complex</b> are the sentences?</p> <p><i>Sentence complexity</i> <i>Text readability</i></p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Sentence/text readability (via textstat)</li> <li>• Sentence complexity</li> </ul>
<p>How is the text vocabulary? Are there <b>unusual words</b>, <b>misspelled words</b>? Are there clear patterns in the <b>first and last few sentences</b> of documents? More nouns, more verbs... Are there <b>repetitions</b>? How rich is an author's vocabulary?</p> <p><i>Vocabulary</i> <i>Unusual words/misspelled words</i> <i>Short/long words, vowel words</i> <i>Vocabulary richness</i> <i>Abstract/concrete vocabulary</i> <i>Objective/subjective writing</i> <i>In which language are texts written in?</i></p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Style analysis (ALL options GUI)</li> </ul> <p>Then tick the Vocabulary analysis option and explore the many options available.</p>
<p><i>Who wrote the text? Man or woman?</i></p>	<p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Who wrote the text? Man or woman? (via Gender guesser)</li> </ul>
<p>Is there <i>dialogue in the text</i>? And who does the talking?</p> <p><i>Is there dialogue in the text?</i></p>	<p>Stanford CoreNLP, uniquely, has a special quote annotator</p>



	<p>that provide answers to these questions.</p> <p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Parsers &amp; annotators (BERT, CoreNLP, spaCy, Stanza)</li> </ul>
<p style="text-align: center;"><b>Visualizing words</b></p>	
<p>How can you get those pretty pictures of words in different colors and sizes?</p> <p><i>Wordclouds</i></p> 	<p>This type of image is called a wordcloud, easily done in the NLP Suite and with a wide range of options.</p> <p><i>Under CORPUS/DOCUMENT analysis tools select:</i></p> <ul style="list-style-type: none"> <li>• Wordclouds (ALL options GUI)</li> </ul>
<p>How can you visualize relations between people or words?</p> <p><i>Sankey plot</i></p> 	<p>This type of plot is called a Sankey plot.</p> <p><i>Under Visualization tools select:</i></p> <ul style="list-style-type: none"> <li>• Sankey flowchart (Plotly) (Open GUI)</li> </ul>
<p>How can you visualize relations between people or words (e.g., Subject-Verb-Object)?</p>	<p>This type of plot is called a network graph and you can display one easily in the NLP Suite (and some NLP tools, for instance SVO, export a network graph automatically).</p> <p><i>Under Visualization tools select:</i></p>

## Network graphs



- Network graphs  
(Gephi) (Open GUI)

Is it possible to visualize gendered names (male or female) in different colors to have an immediate view of where male/female names occur in the text?

## HTML files

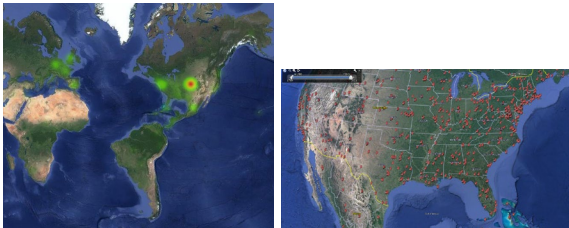
[illegible]

*Under Visualization tools  
select:*

- HTML annotator - dictionary, gender, DBpedia, YAGO, WordNet - (All options GUI)

Can you visualize locations mentioned in texts as geographic maps?

## Pin maps heat maps



*Under  
CORPUS/DOCUMENT  
analysis tools select:*

- Geographic maps:  
From texts to maps

You can also map locations from a list of locations in a csv file:

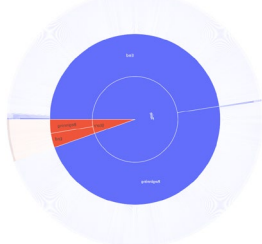
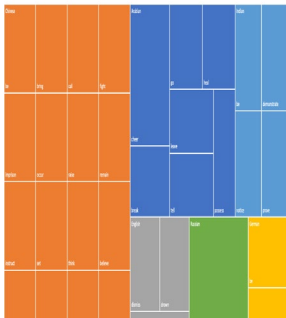

- Geographic maps:  
Google Earth Pro

You will find the same GIS tools under *Visualization tools*.

I have seen presentations with fancy charts like this that move around when you click different parts of the chart. Can it be done easily in the NLP Suite?

## Sunburst pie charts

This type of chart is called sunburst pie chart. It is a powerful interactive visualization tool that you can easily use in the NLP Suite.

	<p><i>Under Visualization tools select:</i></p> <ul style="list-style-type: none"> <li>• Sunburst pie chart (Plotly) (Open GUI)</li> </ul>
<p>How can you produce a map of this kind (from a comparative analysis of folktales across 6 different countries)?</p> <p><i>Treemap</i></p> 	<p>This type of map is called a treemap and you can easily create one in the NLP Suite.</p> <p><i>Under Visualization tools select:</i></p> <ul style="list-style-type: none"> <li>• Treemap (Plotly) (Open GUI)</li> </ul>
<p><i>Animated time plot</i></p>	<p><i>Under Visualization tools.</i></p> <ul style="list-style-type: none"> <li>• Animated time-dependent bar plot (Plotly) (Open GUI)</li> </ul>
 <p><i>Box plots</i></p>	<p><i>Under Visualization tools select.</i></p> <p>Box plots (Open GUI)</p>
<p>Can this type of chart be easily produced to visualize the frequency distribution of certain aspects of texts (e.g., verbs like in this case)?</p> <p><i>Bar charts, pie charts, line charts</i></p>	<p>Every tool in the NLP Suite routinely produces in output a range of plots in either Excel or Plotly, as bar charts or pie charts. You can also access the visualization tools</p> <p><i>Under Visualization tools select:</i></p>

