

Preparing Your Corpus

By Devin Cornell (dcornell@ucsb.edu)

Email me if you have any questions or need help.

[What does a corpus look like?](#)

[How can I download articles from Lexis Nexus?](#)

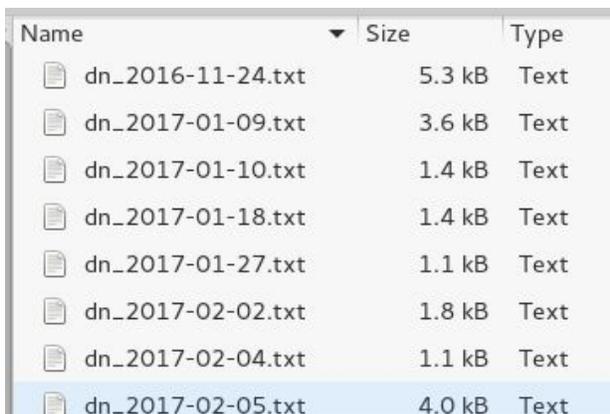
[Can I collect my data any other way?](#)

[Conclusion](#)

What does a corpus look like?

Text analysis is performed on a *corpus*, or a collection of documents. The selection of your corpus is important - it will determine what kinds of questions you can answer and the arguments you can make. For our purposes, a corpus should be prepared as a series of text files. I'll show you how to construct a corpus related to Betsy DeVos for demonstration. Try to imagine how the same instructions can apply to the questions you'd like to answer for yourself.

Our corpora will just be a set of text files in a directory. Text files always end with ".txt", (although your computer may not show the file extension to you). It will still have an icon that looks somewhat like a notepad. Note the filenames are organized by source daily news (dn) and new york times (nyt - files not shown here), then by date yr, month, day. This is important because it will allow us to separate the data by sorting by filename. All of the text files should appear in one folder.



Name	Size	Type
dn_2016-11-24.txt	5.3 kB	Text
dn_2017-01-09.txt	3.6 kB	Text
dn_2017-01-10.txt	1.4 kB	Text
dn_2017-01-18.txt	1.4 kB	Text
dn_2017-01-27.txt	1.1 kB	Text
dn_2017-02-02.txt	1.8 kB	Text
dn_2017-02-04.txt	1.1 kB	Text
dn_2017-02-05.txt	4.0 kB	Text

The context of the text file will simply be the article then. This is a daily news article appears in one text file by itself.

Daily News (New York)

November 24, 2016 Thursday
SPORTS FINAL REPLATE EDITION

TRUMP'S 1ST LADIES * Taps voucher big for Ed. Dept. * S.C. Gov. Haley to represent U.S. at U.N.

BYLINE: BY ADAM EDELMAN NEW YORK DAILY NEWS With Ben Chapman and News Wire Services

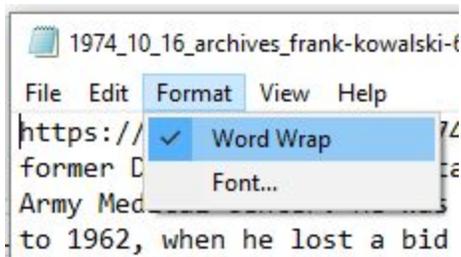
SECTION: NEWS; Pg. 8

LENGTH: 736 words

DONALD TRUMP moved to fill two key cabinet positions Wednesday - naming GOP fund-raiser and school-choice advocate Betsy DeVos education secretary and offering South Carolina Gov. Nikki Haley the position of UN ambassador.

The selections, the first two women to be named to the President-elect's cabinet, raised eyebrows across the political spectrum - with many questioning Haley's credentials and DeVos' policy positions.

Windows users: If you open the text file in notepad, you may need to turn on 'word wrap' for the Betsy Devos articles since the text all appears on one line.



I'll show you how to build this corpus using Lexis Nexus.

How can I download articles from Lexis Nexus?

I'll show you how to build a 'Betsy DeVos' corpus using the Lexis Nexus database. First, find the link to Lexis Nexus on the library page; this is a shortcut: [NexisLexus](#) (this link works if you are on the UCSB campus, not sure about other locations/universities).

Choose to use the current version (will be out of date in 2018 I guess).

Nexis Uni (previously LexisNexis Academic)

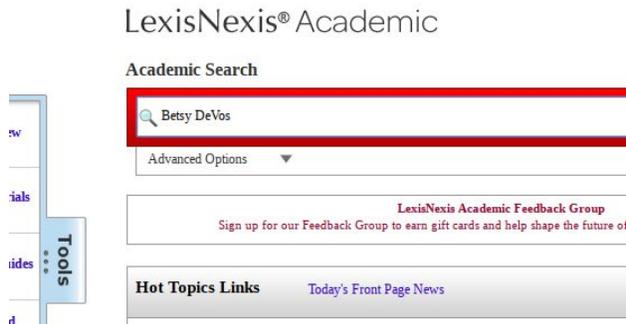
Nexis Uni, which replaces LexisNexis Academic, is available now. UC retains dual acc through December 17, 2017, after which the former will be deactivated across the syst

Go to the new platform: Search Nexis Uni Go to it now	Go to the current platform: Search LexisNexis Academic Go to it now
---	---

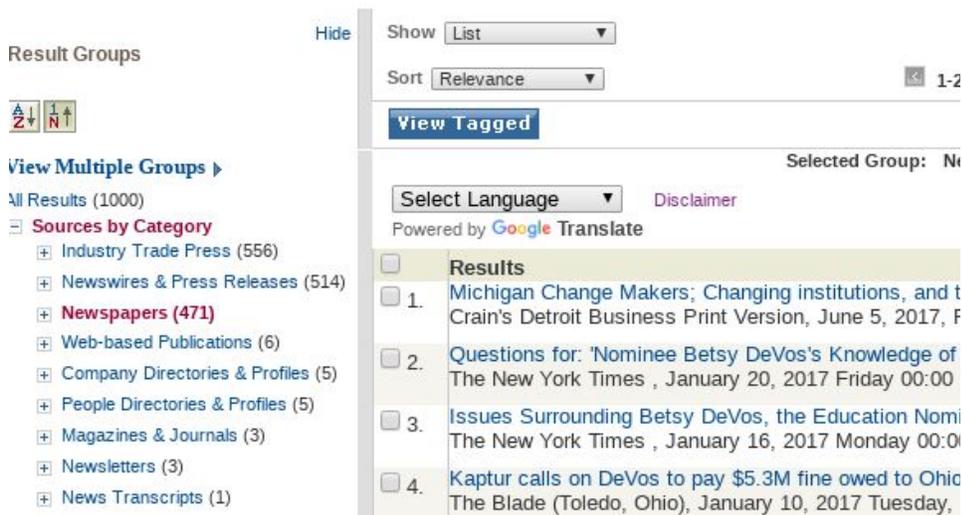
Please contact the Resource Liaison, Harold Colson (hcolson@ucsd.edu), with questio

Access Information

Now start searching. It is up to you to select your query as appropriate for the corpus you want to build. My corpus will be built using “Betsy DeVos” as a search term.



You'll get the following search results. Since I'm interested in Newspapers, I select that on the left-hand column. It appears that there are 471 documents in that category, so it will narrow down the search results.



I'll first build the New York Times part of my corpus. I'll click on the New York Times option under the Newspapers on the left-hand column to narrow the search term. It appears that there are 47 documents returned in the New York Times category.

The screenshot shows a search results interface. On the left, there's a sidebar titled "Result Groups" with a "View Multiple Groups" button. Below it, "All Results (1000)" are listed under "Sources by Category": Industry Trade Press (556), Newswires & Press Releases (514), and Newspapers (471). Under "Newspapers", "The New York Times (47)" is highlighted, with sub-items: McClatchy Tribune non-restricted, Wall Street Journal Abstracts (37), and The Salt Lake Tribune (22). The main area shows search results sorted by "Relevance", displaying items 1-25 of 47. The selected group is "The New York Times". A "View Tagged" button is visible. Below the results, there's a "Select Language" dropdown, a "Disclaimer", and a note "Powered by Google Translate".

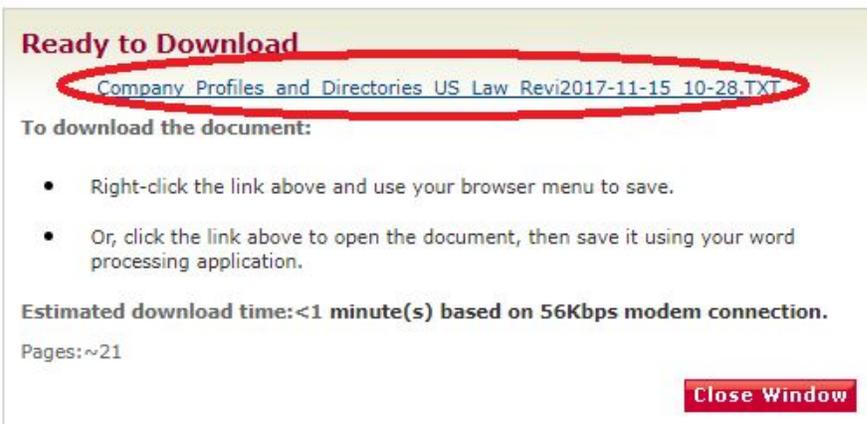
Once you have a list of search results you like, download the entire search query using this button.

This screenshot is similar to the previous one but highlights the download button. The "View Tagged" button is circled in red. Below the search results, there are icons for "Download" (a document with a download arrow) and "Print" (a printer icon). The "Download" button is the one being pointed to by the red circle.

Choose the following download settings for our purposes. Note that it uses text format and I've selected all the documents in this search query (1-47). Once the settings have been configured, hit the download button. It will second to load.

The screenshot shows the "Download Documents" configuration dialog. At the top, it displays the source: "Company Profiles and Directories;US Law Reviews and Journals, Combined;Federal & State Court Case..." and the terms: "(Betsy DeVos)". There are "Download" and "Cancel" buttons. The "Format" is set to "Text". Under "Document View", "Full Document" is selected. Under "Document Range", "Select Items" is chosen, with a text box containing "1-47" and "e.g., 1,3-5,9". Under "Page Options", "Each Document on a New Page" is checked. Under "Font Options", "Times New Roman" is selected, and "Search Terms in Bold Type" is checked. At the bottom, there's a disclaimer: "Download delivery is subject to Terms & Conditions. Please review them. The delivered items will show as activity for the Project ID that initiated the delivery." and "Download" and "Cancel" buttons.

Once the program finishes, you can click the link to download the file. The file will save to somewhere on your computer that you can open - it's just one big text file with all of your articles in it.



Ready to Download

[Company Profiles and Directories US Law Revi2017-11-15 10-28.TXT](#)

To download the document:

- Right-click the link above and use your browser menu to save.
- Or, click the link above to open the document, then save it using your word processing application.

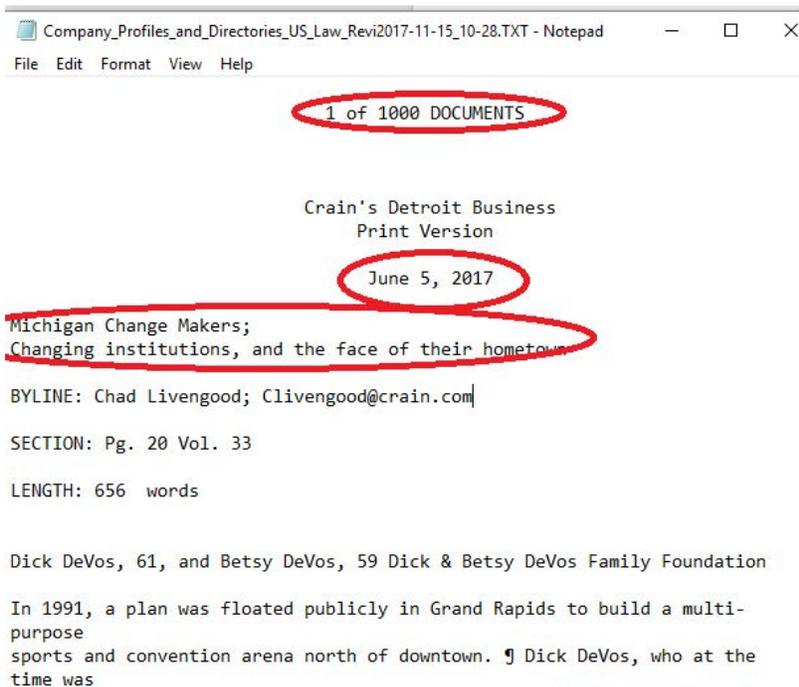
Estimated download time: <1 minute(s) based on 56Kbps modem connection.

Pages: ~21

Close Window

In my case, I wanted two different sub-corpora: one with New York Times, and one with Daily News. In that case, I downloaded the two separate search queries. The next step is to separate the individual articles out of the big file into their own files.

Open the new file and you'll see this.



Company_Profiles_and_Directories_US_Law_Revi2017-11-15_10-28.TXT - Notepad

File Edit Format View Help

1 of 1000 DOCUMENTS

Crain's Detroit Business
Print Version

June 5, 2017

Michigan Change Makers;
Changing institutions, and the face of their hometown

BYLINE: Chad Livengood; Clivengood@crain.com|

SECTION: Pg. 20 Vol. 33

LENGTH: 656 words

Dick DeVos, 61, and Betsy DeVos, 59 Dick & Betsy DeVos Family Foundation

In 1991, a plan was floated publicly in Grand Rapids to build a multi-purpose sports and convention arena north of downtown. ¶ Dick DeVos, who at the time was

All the documents in the query have been added into a single text file, so you'll need to separate them to where one article is one text file. You can do that by copy-pasting the articles into new text files. Every document in this text file starts with a "x of 1000 DOCUMENTS", so to get all of

the document content, you should just copy-paste everything in between the headers for consecutive documents.

```
2 of 1000 DOCUMENTS

The New York Times
January 20, 2017 Friday 00:00 EST
Questions for: 'Nominee Betsy DeVos's Knowledge of Education Basics
Criticism';
Article of the Day
BYLINE: CAROLINE CROSSON GILPIN
SECTION: LEARNING
LENGTH: 387 words

HIGHLIGHT: What do you think about the prospect of Ms. DeVos heading
country's Department of Education?

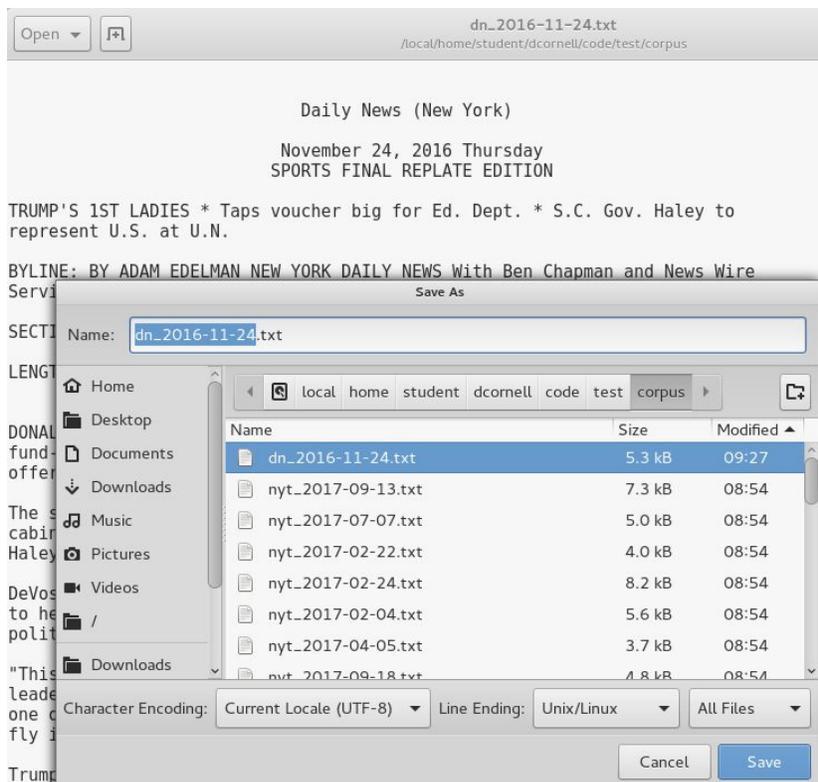
Article: Nominee Betsy DeVos's Knowledge of Education Basics Is Open
Criticism

Before Reading

The United States Department of Education outlines its mission as fo
The mission of the Department of Education is to ensure equal access
education and to promote educational excellence throughout the natio
```

Open up a blank text file in notepad or whatever text editor you use. If you open those programs, it should start with a blank file that you can paste into. Paste in your article, then save the article with an appropriate filename.

The best filenames include the source name/abbreviation (newspaper), a date (for chronological sorting), and the title of the article (or some other unique identifier). The above text file is used as an example. I didn't use any unique document identifier, but you may want to - especially if you have two documents from the same day.



The naming of the file is important - it is called 'metadata'. It is data about the data we want to work with. When we go to use the topic modeling software, it will be especially important because it will help us sort through the documents for close reading or quantitative analysis.

Can I collect my data any other way?

Yes, but it all needs to end up in the same format: a series of text files with appropriate names. For instance, instead of copying from the Lexis Nexus download page, you can copy-paste an article from the web in the same way. If you copy into a program like 'notepad', it will just ignore the image (this is what we want).

Isla Vista Landlord Charged With Criminal Complaint After Yelling Homophobic Slurs at CSD President

November 15, 2017 at 12:34 am by Jose Ochoa

Isla Vista landlord James Gelb will be charged with disturbing the peace after he was filmed yelling homophobic slurs at CSD President Ethan Bertrand.

Santa Barbara County District Attorney Joyce Dudley said she would file charges against Gelb. Dudley said the charge is a response to a Nov. 11 incident in which Gelb used offensive language in order to likely provoke Bertrand.

Disturbing the peace can result in imprisonment in county jail for up to 90 days or a fine of up to \$400, according to California law.

Gelb is scheduled to appear in court for an arraignment on Dec. 4. He told the Nexus on Tuesday that he has hired an attorney to "aggressively defend" himself from the "unfounded" accusations.

- Copy Ctrl+C
- Search Google for "Isla Vista Landlord Charged With Criminal..."
- Print... Ctrl+P
- AdBlock
- Google Translate
- Inspect Ctrl+Shift+I

Tweets by @dailynexus

Daily Nexus @dailynexus
UCSBATL: Police determined that the threat was "false and unfounded" after searching the area, according to UCSD Sgt. Dan Wilcox. [dailynexus.com/2017/11/14/news...](#)

Conclusion

In the end, the most important thing is to have a list of text files you can analyze, with filenames that are useful for your analysis (i.e. convenient names for sorting). It is recommended that you have at least around 100 text articles to perform topic modeling or sentiment analysis (although, feel free to use many more - the software can use as many as you can make).

After you have finished your corpus, upload it to Gauchospace and we'll run the topic modeling software on it. You can read the 'instructions overview' to understand how to analyze the data. Get the corpus collected as soon as possible so you can get started with the analysis.

If you have any technical questions, feel free to email dcornell@ucsb.edu.