

User Manual

sec-10-k-scraper

Creators: Ed Carl, Chloe Dorward, Seth Markarian, Sena Yevenyo

An up-to-date version of this document can always be found [on Google Docs](#).

A video demonstration of the app can be found [here](#).

Table of Contents

[Table of Contents](#)

[Summary](#)

[Purpose of Application](#)

[The App](#)

[Installing the App](#)

[Windows](#)

[Mac](#)

[Linux](#)

[Running The App](#)

[Step 1: Selecting Filings](#)

[Option 1: Searching](#)

[Searching](#)

[Selecting Filings](#)

[Option 2: Uploading](#)

[File Format](#)

[Uploading the File](#)

[Step 2: Viewing Queue](#)

[Step 3: Downloading Files](#)

[Step 4: Finding/Using Your Files](#)

[File Structure](#)

[Excel File](#)

[Errors](#)

[Error Dialog](#)

[Search Error Dialog](#)

[Queue Error Dialog Box](#)

[Developer Tools](#)

Summary

Purpose of Application

This app is a local desktop application for Windows, MacOS, and Linux that obtains companies' 10-K, 10-Q, and 10-F SEC filings, downloads the filings as HTML files, extracts the text of key sections of those documents, analyzes the text for named entities, and gathers that information in a Microsoft Excel spreadsheet. Users can either search for filings based on company names/ CIK numbers, filing type, and date range, or upload a file including similar parameters for batch uploading.

Using a desktop application rather than a web service minimizes the technical and financial requirements for FracTracker to operate it. The application is entirely self-contained on users' computers and there will be no need to spend time on configuring it or money on hosting a server. A user only needs to have an internet connection in order to run the application.

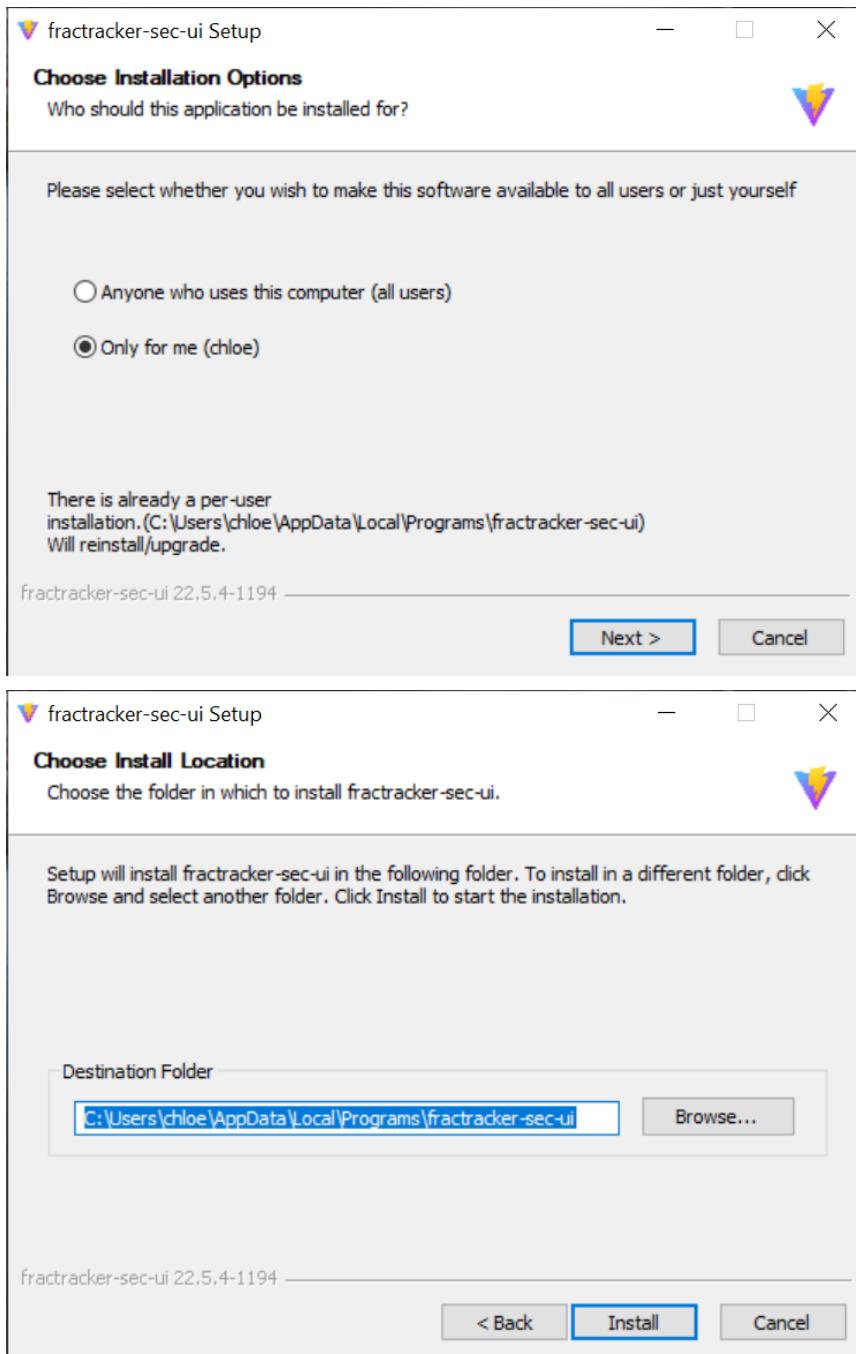
The App

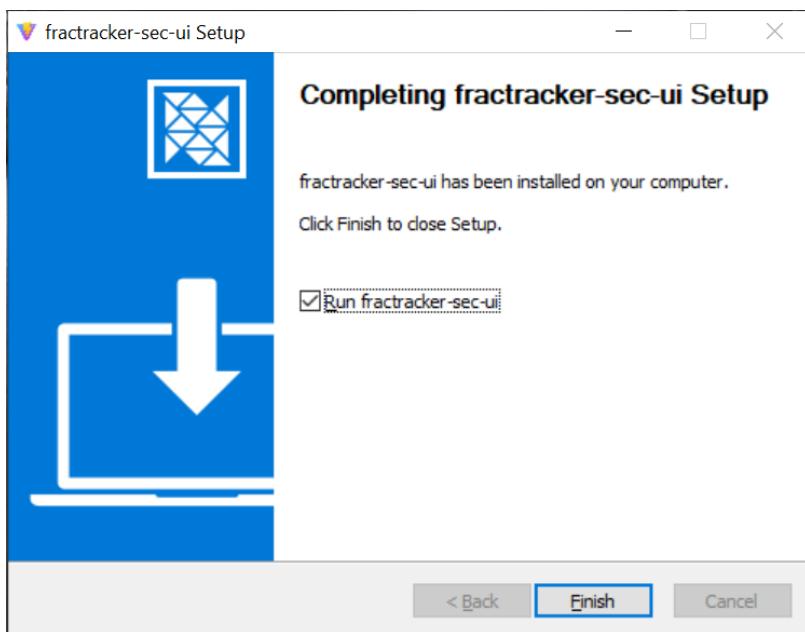
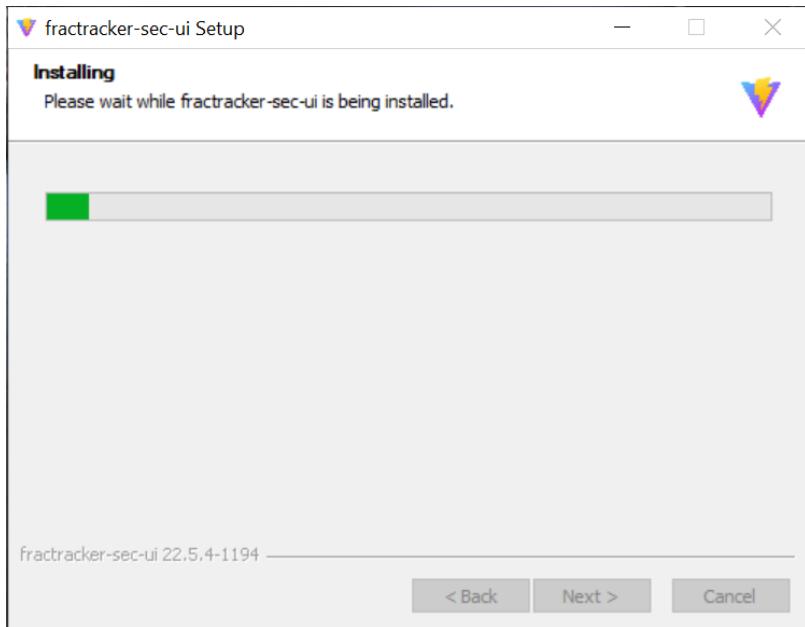
Installing the App

Releases of the application installer can be found under the releases page of the project's github repository. github.com/edwardmcarl/sec-10-k-scraper/releases

Windows

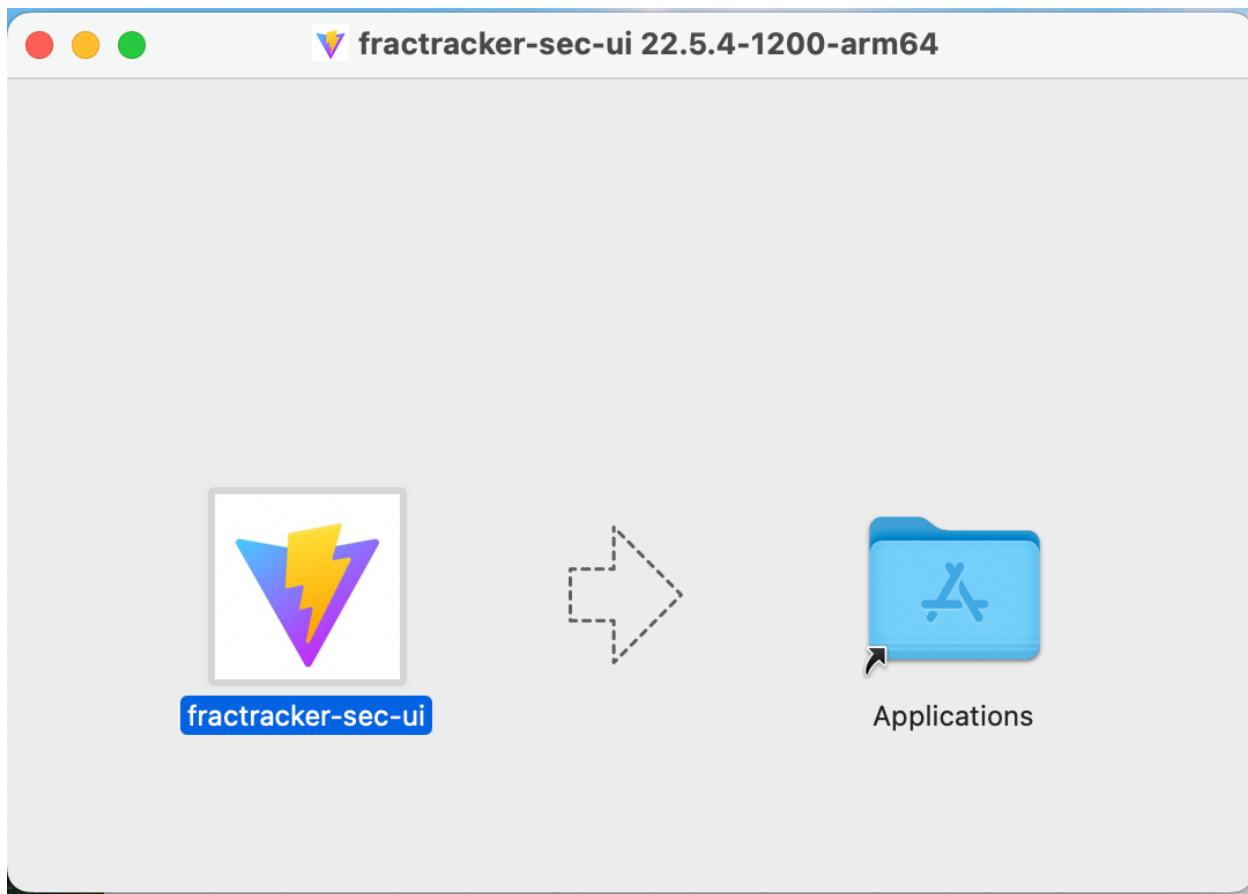
Find the installer in your file explorer (fractracker-sec-ui Setup 22.5.4-1194.exe), and run it. Navigate through the prompts to install.





An icon will appear on your desktop.

Mac



Simply mount the installer .dmg and drag the app into the Applications folder.

Linux

Just double-click the installer AppImage.

Running The App

Step 1: Selecting Filings

Option 1: Searching

Searching

The first option is to select an entity, and search for filings based on a date range.



SEC EDGAR 10-K Information Access

ford

CIK0000038264 | Forward Industries, Inc. (FORD)

CIK0000037996 | FORD MOTOR CO (F, F-PC, F-PB)

CIK0001586097 | Ford Tim

CIK0001176401 | FORD JOE T

CIK0001216553 | FORD BETH

CIK0001511094 | Ford Adam

CIK0001390145 | Ford Jill

CIK0001606824 | Ford Mark

CIK0001392324 | Ford Eric

CIK0001143133 | FORD MELBA

Click on the search bar, and begin to input the appropriate entity name or CIK number. A list of predicted entities will appear in a list below the bar.



SEC EDGAR 10-K Information Access

FORD MOTOR CO (F, F-PC, F-PB)

Start Date: 5/2/2021

End Date: 5/2/2022

Form Type:

Read from file (.txt):

No file chosen

Entity Name	CIK Number	Form Type	Filing Date	Link to Document	Extract Info?
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	https://sec.gov/Archives/edgar/data/37996/000003799622000013/f-20211231.htm	<input type="button" value="Add to Queue"/>

Select an entity by clicking on it (for example, FORD MOTOR CO). This will populate the results table with the default search parameters - 10-K filings from the past year.

Start Date: 5/2/2021 End Date: 5/2/2022

Read from

Entity Name	CIK Number	Form Type	Filing Date	Link to Document
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	https://sec.gov/Archives/edgar/data/37996/000003799622000013/f-20211231.htm

Show Queue

May 2021

MON	TUE	WED	THU	FRI	SAT	SUN
26	27	28	29	30	1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31	1	2	3	4	5	6

You can adjust your date range by clicking on the date selector, and changing the date accordingly.

Start Date: 5/2/2022 End Date: 5/2/2022

Read from

MON	TUE	WED	THU	FRI	SAT	SUN
29	30	1	2	3	4	5
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31	1	2

Entity Name CIK Number Form Type Filing Date

You can also type in the date rather than using the provided selector.

Form Type:

10-K

10-Q

20-F

You can change the selected form type by clicking on the gray dropdown on the right hand side, and clicking on the appropriate form.

Entity Name	CIK Number	Form Type	Filing Date	Link to Document	Extract Info?
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	https://sec.gov/Archives/edgar/data/37996/000003799622000013/f-20211231.htm	<input type="button" value="Add to Queue"/>
FORD MOTOR CO	CIK0000037996	10-K	2021-02-05	https://sec.gov/Archives/edgar/data/37996/000003799621000012/f-20201231.htm	<input type="button" value="Add to Queue"/>
FORD MOTOR CO	CIK0000037996	10-K	2020-02-05	https://sec.gov/Archives/edgar/data/37996/000003799620000010/f1231201910-k.htm	<input type="button" value="Add to Queue"/>

Clicking on the blue Search button will repopulate the results table based on your selected parameters. Changing the date, form type, or search components also repopulates the results.

Selecting Filings

Entity Name	CIK Number	Form Type	Filing Date	Link to Document	Extract Info?
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	https://sec.gov/Archives/edgar/data/37996/000003799622000013/f-20211231.htm	<input type="button" value="Add to Queue"/>
FORD MOTOR CO	CIK0000037996	10-K	2021-02-05	https://sec.gov/Archives/edgar/data/37996/000003799621000012/f-20201231.htm	<input type="button" value="Add to Queue"/>
FORD MOTOR CO	CIK0000037996	10-K	2020-02-05	https://sec.gov/Archives/edgar/data/37996/000003799620000010/f1231201910-k.htm	<input type="button" value="Add to Queue"/>

Individual filings from different years can be selected by clicking on the Add to Queue button.

Entity Name	CIK Number	Form Type	Filing Date	Link to Document	Extract Info?
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	https://sec.gov/Archives/edgar/data/37996/000003799622000013/f-20211231.htm	<input type="button" value="Remove from Queue"/>
FORD MOTOR CO	CIK0000037996	10-K	2021-02-05	https://sec.gov/Archives/edgar/data/37996/000003799621000012/f-20201231.htm	<input type="button" value="Remove from Queue"/>
FORD MOTOR CO	CIK0000037996	10-K	2020-02-05	https://sec.gov/Archives/edgar/data/37996/000003799620000010/f1231201910-k.htm	<input type="button" value="Remove from Queue"/>

Individual filings can be removed by clicking on the Remove from Queue button.
Repeat the process outlined above to add filings from different entities to the same Queue.

Option 2: Uploading

The second option is to upload a text file with different CIK numbers and appropriate date ranges. This option will add every filing within a CIK's date range to the Queue.

File Format

The file must be a .txt file. Your computer will have a program to create text files, for example, Notepad,TextEdit, etc.

Each line represents a new entity, and should contain the CIK number, the start of the date range, and the end of the date range:

CIKXXXXXXXXXX YYYY-MM-DD YYYY-MM-DD where X is 0-9 inclusive.

Uploading the File

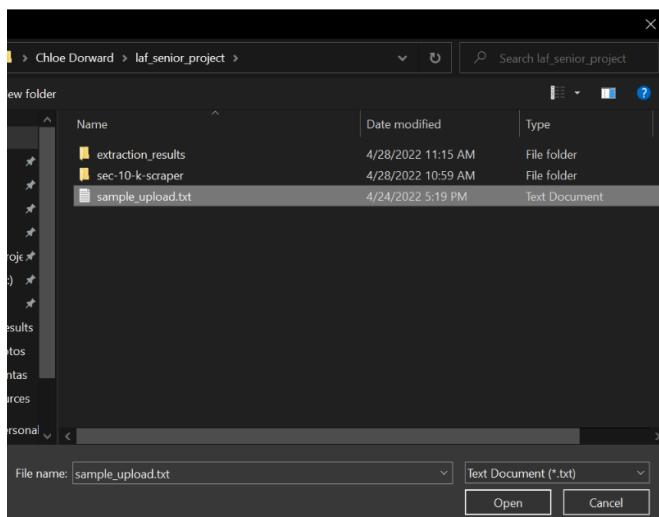
Read from file (.txt): [\(i\)](#)

[Choose File](#) No file chosen

No file chosen

[Search](#)

Click on the Choose File button.



Navigate through your file system to select your file.

Read from file (.txt): [\(i\)](#)

[Choose File](#) sample_upload.txt

[Search](#)

The uploaded file name will appear next to the Choose File button.

This option can be used in combination with Option 1. Any results that have already been selected through the File Upload option will appear as such in the search results table.

The screenshot shows the SEC EDGAR 10-K Information Access interface. A modal window titled "File Upload Setup" is open. It contains fields for "Start Date" (5/2/2019), "End Date" (5/2/2022), and "Form Type" (10-K). Below these are instructions: "[CIK] [START_DATE] [END_DATE]" and "[CIK] [START_DATE] [END_DATE]". A note states "* Dates are in ISO format (YYYY-MM-DD)". There is also a "Choose File" button with the path "sample_upload.txt" and a "Search" button.

If you forget the format for file uploading, it can be found by clicking on the i button.

Step 2: Viewing Queue

The screenshot shows the SEC EDGAR 10-K Information Access interface. A modal window titled "Viewing Queue" is open. It contains fields for "Start Date" (5/2/2019), "End Date" (5/2/2022), and "Form Type" (10-K). Below these are instructions: "[CIK] [START_DATE] [END_DATE]" and "[CIK] [START_DATE] [END_DATE]". A note states "* Dates are in ISO format (YYYY-MM-DD)". There is also a "Choose File" button with the path "sample_upload.txt" and a "Search" button.

Entity Name	CIK Number	Form Type	Filing Date	Link to Document	Extract Info?
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	https://sec.gov/Archives/edgar/data/37996/000003799622000013/f-20211231.htm	<button>Remove from Queue</button>
FORD MOTOR CO	CIK0000037996	10-K	2021-02-05	https://sec.gov/Archives/edgar/data/37996/000003799621000012/f-20201231.htm	<button>Remove from Queue</button>
FORD MOTOR CO	CIK0000037996	10-K	2020-02-05	https://sec.gov/Archives/edgar/data/37996/000003799620000010/f1231201910-k.htm	<button>Remove from Queue</button>

[Show Queue](#)

View the Queue by clicking on the blue Show Queue button on the bottom of the page.

The screenshot shows the Fractracker Alliance SEC EDGAR 10-K Information Access interface. The main window displays search parameters (Start Date: 5/2/2019, End Date: 5/2/2022) and a table of results for Ford Motor Co. The Queue tab is open on the right, showing a list of filings with 'Remove From Queue' buttons.

Entity Name	CIK Number	Form Type	Filing Date	Link to Document
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	https://sec.gov/Archives/edgar/data/37996/0
FORD MOTOR CO	CIK0000037996	10-K	2021-02-05	https://sec.gov/Archives/edgar/data/37996/0
FORD MOTOR CO	CIK0000037996	10-K	2020-02-05	https://sec.gov/Archives/edgar/data/37996/0

Show Queue

Entity Name	CIK Number	Form Type	Filing Date	
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	Remove From Queue
FORD MOTOR CO	CIK0000037996	10-K	2021-02-05	Remove From Queue
FORD MOTOR CO	CIK0000037996	10-K	2020-02-05	Remove From Queue
HERSHEY CO	CIK0000047111	10-K	2020-02-20	Remove From Queue
HERSHEY CO	CIK0000047111	10-K	2019-02-22	Remove From Queue
CONTINENTAL RESOURCES, INC	CIK0000732834	10-K	2021-02-16	Remove From Queue
CONTINENTAL RESOURCES, INC	CIK0000732834	10-K	2020-02-26	Remove From Queue

This will slide open the Queue tab.

Entity Name	CIK Number	Form Type	Filing Date	
FORD MOTOR CO	CIK0000037996	10-K	2022-02-04	Remove From Queue

From here, you can remove filings from the queue.

Step 3: Downloading Files

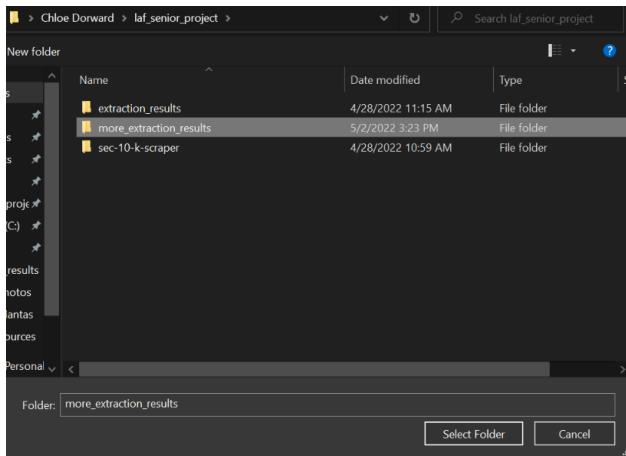
Queue

Choose Path for Download
C:\Users\chloe\Desktop

Apply Named Entity Recognition to Queue

Extract & Download

At the top of the Queue tab, click the Choose Path button. This allows you to select a place to save your extracted files. The default path is your Desktop folder.



Navigate through the file dialog, and select a folder to save the files to.

Queue

Choose Path for Download
C:\Users\chloe\laf_senior_project\more_extraction_results
 Apply Named Entity Recognition to Queue
Extract & Download

Once a location is selected, the name of the file path will appear below the Choose File button.

Queue

Choose Path for Download
C:\Users\chloe\laf_senior_project\more_extraction_results
 Apply Named Entity Recognition to Queue
Extract & Download

Next, you have the option to either apply Named Entity Recognition during the file extraction process, or not. You can select the checkbox to do this.

Queue

Choose Path for Download
C:\Users\chloe\laf_senior_project\more_extraction_results
 Apply Named Entity Recognition to Queue
Extract & Download

Finally, press the Extract & Download button.

The screenshot shows the Fractracker Alliance SEC EDGAR 10-K Information Access interface. On the left, there's a search form with date filters (Start Date: 5/2/2019, End Date: 5/2/2022), a file upload section (Choose File: sample_upload.txt), and a table of search results for Ford Motor Co. On the right, a 'Queue' window is open, showing a table of filings with a 'Remove From Queue' button next to each entry.

While the system is extracting the files, it will indicate that it is loading. The buttons next to the selected filings will gray out, meaning that they are processing. You cannot deselect an entity once the extraction process has started.

Once the extraction process is complete, the loading spinner will disappear, and the selected filings will no longer be grayed out.

Step 4: Finding/Using Your Files

File Structure

Files will be downloaded to the folder that you selected in Step 3.

The screenshot shows a file explorer window displaying the directory structure. The root folder contains three subfolders: 'CONTINENTAL RESOURCES, INC', 'FORD MOTOR CO', and 'HERSHEY CO'. There is also a file named 'summary.xlsx'.

Within your selected folder, there will be a folder for each entity, and an excel spreadsheet.

The screenshot shows a file explorer window displaying the contents of the 'CONTINENTAL RESOURCES, INC' folder. It contains two files: '10-K_2020-02-26.htm' and '10-K_2021-02-16.htm'.

Each entity folder will contain htm files for each filing. The name of the file will be the filing type and the filing date.

Excel File

	A	B	C	D	E	F	G	H
1	Company Name	EIN	HQ Address of Incorporation	Business Type	Risk Factor 1	Risk Factor 2	Risk Factor 3	Risk Factor 4
2	FORD MO	380545190	{'street1': 'DE	See origin: []	See origin: ['Autc			
3	FORD MO	380545190	{'street1': 'DE	See origin: []	See origin: ['Autc			
4	FORD MO	380545190	{'street1': 'DE	See origin: []	See origin: ['Autc			
5	HERSHEY	(230691590	{'street1': 'DE	>Item 1.BI []	>Item 1A.II['Amp			
6	HERSHEY	(230691590	{'street1': 'DE	>Item 1.BI []	>Item 1A.II['Amp			
7	CONTINEN	730767545	{'street1': 'OK	See origin: []	See origin: ['App			
8	CONTINEN	730767545	{'street1': 'OK	See origin: []	See origin: ['Audi			

Within the summary.xlsx excel file, each new filing will add a new row to the spreadsheet. If you choose to download your files to a folder which already has the summary.xlsx file, all additional filings will be appended to the end of the spreadsheet; our tool does not look through to check if there are duplicates.

The Columns in the spreadsheet are:

- Company Name
 - EIN
 - HQ Address
 - State of Incorporation
 - 1. Business
 - NER 1. Business
 - 1A. Risk Factors
 - NER 1A. Risk Factors
 - 2. Properties
 - NER 2. Properties
 - 3. Legal Proceedings
 - NER 3. Legal Proceedings
 - 6. [Reserved]
 - NER 6. [Reserved]
 - 7. Financial Condition
 - NER 7. Financial Condition
 - 7A. Market Risk
 - NER 7A. Market Risk
 - 10. Corporate Governance
 - NER 10. Corporate Governance
 - 12. Security Ownership and Stockholder Matters
 - NER 12. Security Ownership and Stockholder Matters
 - 13. Certain Relationships
 - NER 13. Certain Relationships

Cells pertaining to entity information should always be filled in.

Cells pertaining to subsections of the filing will contain the text within that section if possible. If the number of characters in a section happens to be too long for a single Excel cell (32,767 characters) the statement 'See original document; Section too long for Excel' will be placed in the cell.

R	S	T	U	V	W	X	Y	Z
A. Market orporate Gship and Stership andain Relatiertain Relationships ['Automotive', 'Board of Directors', 'CDS', 'Controller', 'Ford Credit', 'Ford Credit's', 'GRMC', 'Treasurer', 'Treasurer's Office', 'the Asset-Liability Committee', 'the Audit Committee', 'the Global Risk Management Committee', 'the London Interbank Offered Rate']								

Cells pertaining to Named Entity Recognition will contain a bracketed, comma-separated list of named entities. If you did not select the Named Entity Recognition checkbox in Step 3, cells under the NER columns will be blank.

Errors

Error Dialog

There are two error dialogs in the application:

- Search Error Dialog
- Queue Error Dialog

Search Error Dialog

The search error dialog box is found in the main page of the application. This light red box with error text in deep red pops up during search for companies' documents that are to be added to the queue for extraction. An image of the error box is found below:

The screenshot shows the FRACTRACKER SEC EDGAR 10-K Information Access search interface. At the top, there is a navigation bar with links for Home, Help, Log Out, and a user profile icon. Below the navigation bar is a search bar with the placeholder "FORD MOTOR CO (F, F-PC, F-PB)". Underneath the search bar are several input fields: "Start Date" set to 5/4/2021, "End Date" set to 5/4/2023, "Form Type" set to 10-K, and a "Read from file (.txt)" input field containing "sample_upload.txt". There is also a "Choose File" button and a "Search" button. A prominent error message "End date cannot be later than today's date" is displayed in a red box. At the bottom, there is a table with columns for Entity Name, CIK Number, Form Type, Filing Date, Link to Document, and Extract Info?.

Entity Name	CIK Number	Form Type	Filing Date	Link to Document	Extract Info?

[Show Queue](#)

Explanation of the errors of the application are listed below:

- **The SEC EDGAR server could not process the request**
 - This means that the SEC server could not process your search query. The server might be down, blocking all requests or there have been too many requests to the server (occasionally happens). Simply try again.
- **The application failed to reach the server**
 - This means that there might be no internet connection. Check your internet connection.
- **CIK number input not in correct format**
 - Happens often during batch processing. This means that the CIK number inputted is not in the structure stated under [Option 2: File Format under Step 1: Selecting Files](#).
- **Start date input not in ISO format**
 - Happens often during batch processing. The start date given in the file is not in ISO 8601 format. ISO format for the dates are: YYYY-MM-DD
- **End date input not in ISO format**
 - Happens often during batch processing. The end date given in the file is not in ISO 8601 format. ISO format for the dates are: YYYY-MM-DD
- **Start date cannot be earlier than 1994-01-01**
 - This is the earliest date of records in the SEC EDGAR database server
- **End date cannot be later than today's date**
 - A user cannot query documents from future dates
- **Start date cannot be later than end date**
 - This error means exactly what it is. Please adjust the start state to be earlier than the end date
- **Something occurred when decompressing and/or decoding response from SEC EDGAR server**
 - Rarely would this happen. If it does there is a change in the way the server is compressing data. Please reach out to Lafayette College Computer Science Department if this happens.
- **The CIK number input does not exist in SEC EDGAR database**
 - Happens often during batch processing. The CIK number is in a correct format but does not exist in the SEC server
- **The form type inputted is not supported by this application**
 - Rarely happens. This means that the form input is no longer supported by the SEC server for querying.

Note that when you submit a text file for batch processing to add companies filing to the queue, if there are errors with processing the text file, any error that occurs when processing each line of the text file will be in the search dialog box. An example of the error would be like:

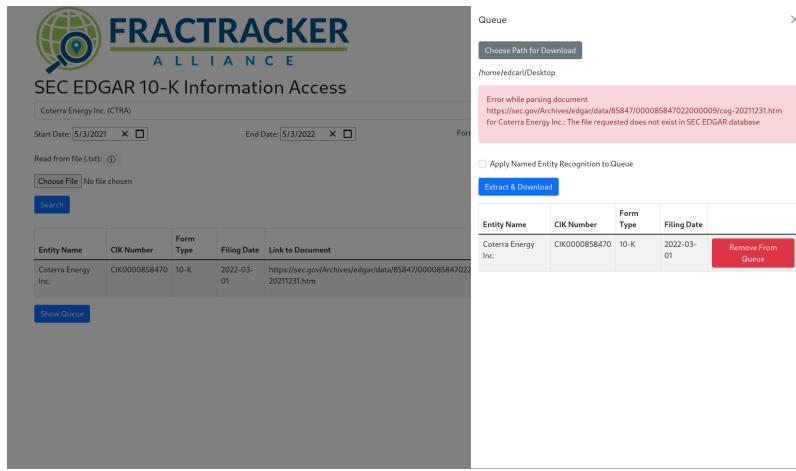
Line 14: CIK number input not in correct format*

*given that Line 14 in your text file for batch processing/file upload does not have the correct CIK number

If multiple errors are encountered, they are all placed in the error box.

Queue Error Dialog Box

The queue error dialog box is found in the side draw of the application which is shown when the user clicks the 'Show Queue' button. This error message pops up below the text of the selected path as shown below:



A list of possible errors are in the possible structure:

Error while {stage where error occurred} document {document URL} for {company}: {Error message}

Error sections explained below:

{stage where error occurred}

The stage where error has occurred possible messages are:

- parsing
 - This means that the error happened during parsing of the document for the itemized sections
- downloading
 - This means that the error occurred when downloading data from the SEC server
- applying NER to
 - This means that the error occurred when applying named entity recognition to the extracted itemized fields
- adding spreadsheet row for
 - This means the error occurred when appending the new data to the output Excel file

{document URL}

This is just the document URL being handled for which the error occurred.

{company}

This is the company's name whose document is currently handled.

{Error message}

Why was the error encountered. The possible messages are:

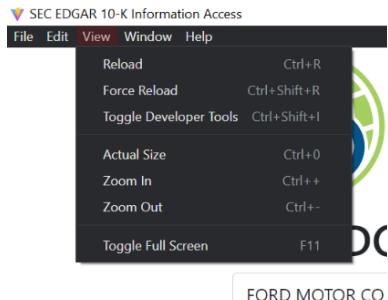
- **Parsing for this document is not supported**
 - Rarely happens. This means that the document is not an HTML (htm. or .html) file and thus cannot be parsed.
- **The SEC EDGAR server could not process the request**
 - This means that the SEC server could not process your document query. The server might be down, blocking all requests or there have been too many requests to the server (occasionally happens). Simply try again.
- **The application failed to reach the server**
 - This means that there might be no internet connection. Check your internet connection.
- **Something occurred when decompressing and/or decoding response from SEC EDGAR server**
 - Rarely would this happen. If it does there is a change in the way the server is compressing data. Please reach out to Lafayette College Computer Science Department if this happens.
- **The file requested does not exist in SEC EDGAR database**
 - Occasionally happens. This means that the file does not exist in the SEC EDGAR database.

A possible message then is:

*Error while downloading document <https://example.com/file-does-not-exist.html> for Example Inc:
The SEC EDGAR server could not process the request.*

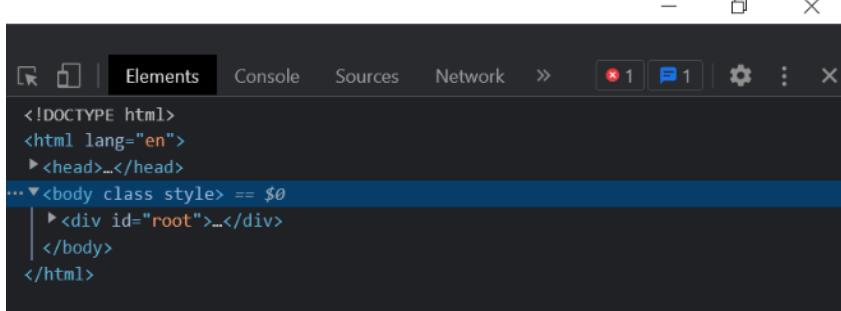
Developer Tools

If you are experiencing errors and cannot tell why, you may be able to find out what is happening by viewing the developer tools.



On the top menu, click on the View tab.

Click the Toggle Developer Tools option.



Click the Console tab.

SEC EDGAR 10-K Information Access

File Edit View Window Help

 FRACTRACKER ALLIANCE

SEC EDGAR 10-K Information Access

FORD MOTOR CO (F, F-PC, F-PB)

Start Date: 4/20/2022

End Date: 4/23/2022

Form Type:

Read from file (.txt): sample_upload.txt

Entity Name	CIK Number	Form Type	Filing Date	Link to Document	Extract Info?
					<input type="button" value="Show Queue"/>

Console Sources Network

Filter

```
> Array(10)
searched
> Array(3)
CIX0000037996
2020-01-01
2022-03-28
> object
> Array(3)
> L
> L
> L
2019-01-01
CIX0000047111
2020-12-31
> object
> Array(2)
> L
> L
CIX00000732834
2020-01-01
2021-12-31
> object
> Array(2)
> L
> L
C:\users\chloe\laf_senior_project\extraction_results\react-dom.production.min.js:101
● uncaught TypeError: Cannot read properties 'react-dom.production' of undefined
at Object.g_ (react-dom.production.min.js:52:312)
at w (react-dom.production.min.js:52:311)
at _ (react-dom.production.min.js:53:35)
at t (react-dom.production.min.js:100:68)
at tv (react-dom.production.min.js:101:380)
at Pd (react-dom.production.min.js:102:189)
at t (react-dom.production.min.js:292:189)
at v_ (react-dom.production.min.js:50:52)
at IV (react-dom.production.min.js:105:469)
c:\users\chloe\laf_senior_project\extraction_results\react-dom.production.min.js:359
NER: true
> |
```

This will display statements that get written out to the application's console, including uncaught error messages.