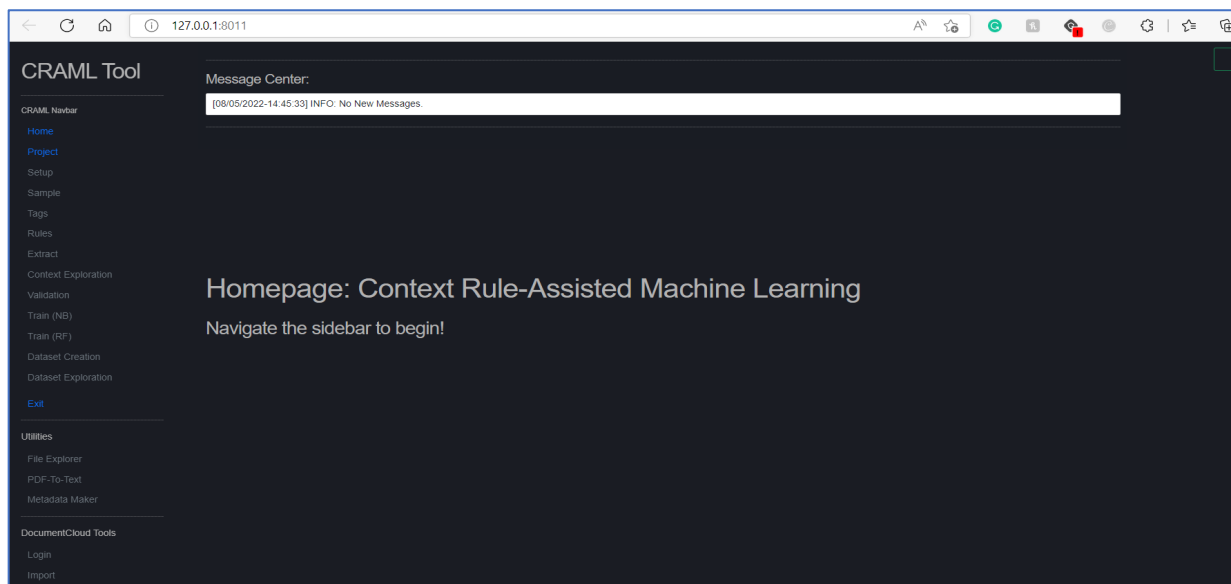**A Step-by-Step Guide to Recreate the No Poach analysis[1]**

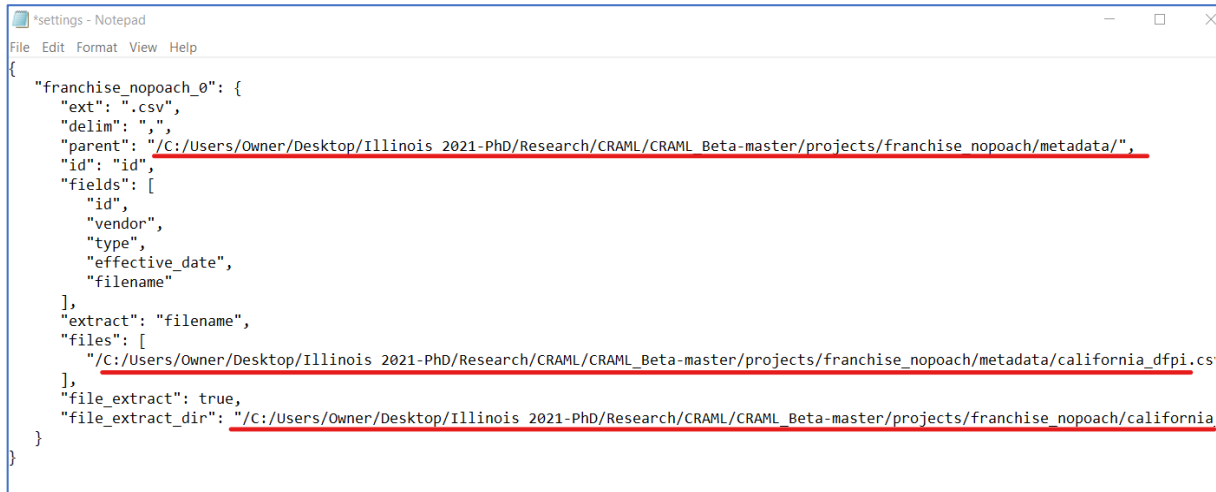1.  Start CRAML_Tool.py with Python. A browser window will appear.



2.  Download the unstructured text data and metadata to a folder.

3.  Download and unzip the intermediate project files.

    a)  In the folder CRAML is installed in, a "projects" folder contains all research projects.

b)  Download the intermediate files and unzip folders that contain the metadata, keywords, and the settings files of the project "franchise_nopoach." Once downloaded, extract the files, and save them in the "projects" folder created previously.

c)  In the settings.json file in each project directory, update the path to point it to the directory that has the cleaned data and the file that contains the metadata. For example, see below:



## Additional Guidance on Working with CRAML

Users can work with their own data inside CRAML.

## Interacting with CRAML

a)  New directories created in the "projects" folder outside of CRAML will appear as projects in CRAML. Project directories deleted outside of CRAML will be deleted in CRAML.

b)  Likewise, creating a new project inside CRAML creates a directory in the projects folder. As users perform analysis in CRAML, new folders and files are created in each project folder to save the output and enable users to interactively engage with the output.

## Setting Up a project on CRAML:

Setup is a challenging part of working in CRAML because text comes in many varieties and there's no one standard way of having a massive text corpus.

Inputs must be in a form that CRAML will accept. CRAML is designed to be text format agnostic, and accepts CSV, JSON, XML and other formats, but for unusual cases, users need some programming experience to input data into CRAML.

"Extract fields contains filename?" refers to whether the metadata refers to specific text files.

"Fields to keep:" refers to whether output datasets should include metadata fields.



## Sampling Data

Sampling selects a percentage of documents and stores the results in sample.csv for further analysis.
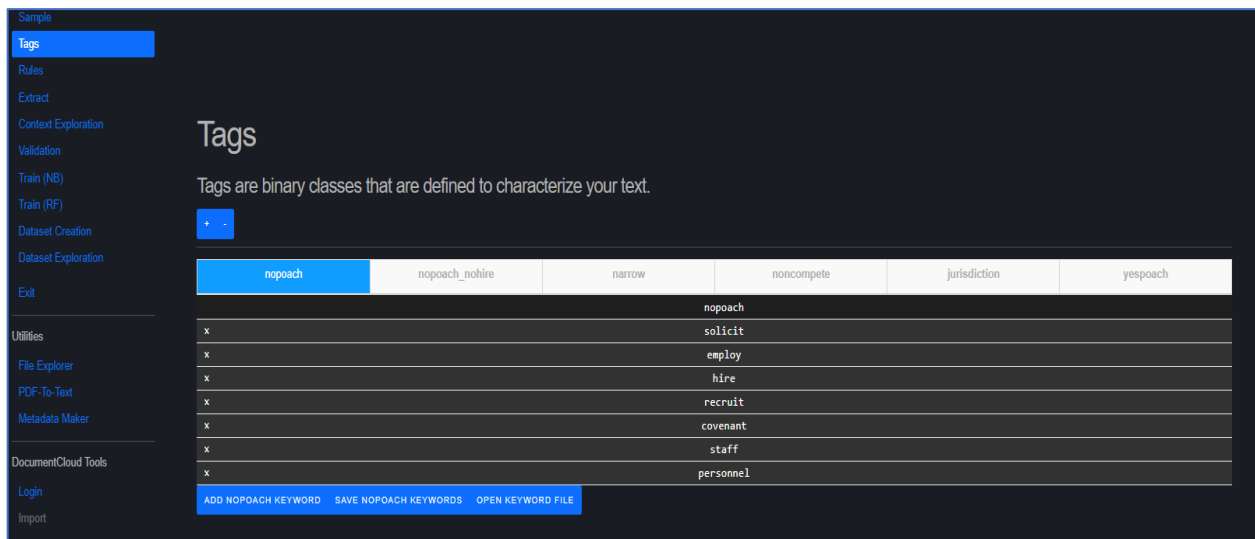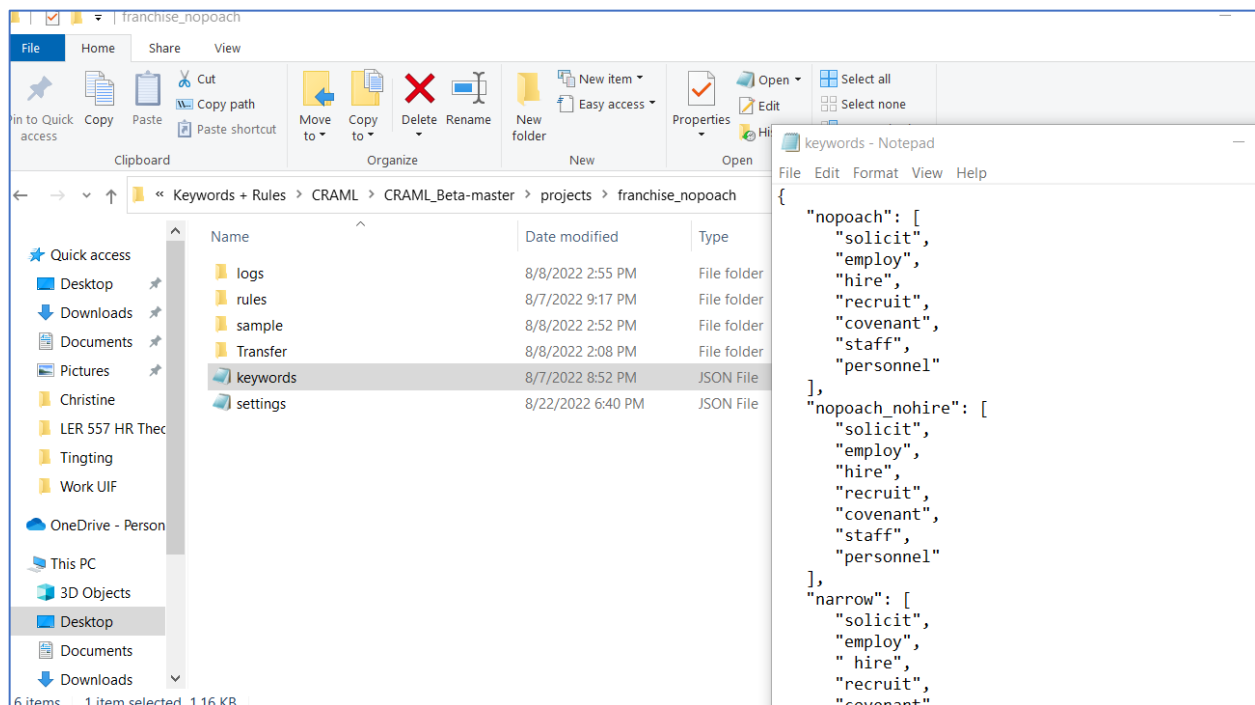
## Tags and Keywords
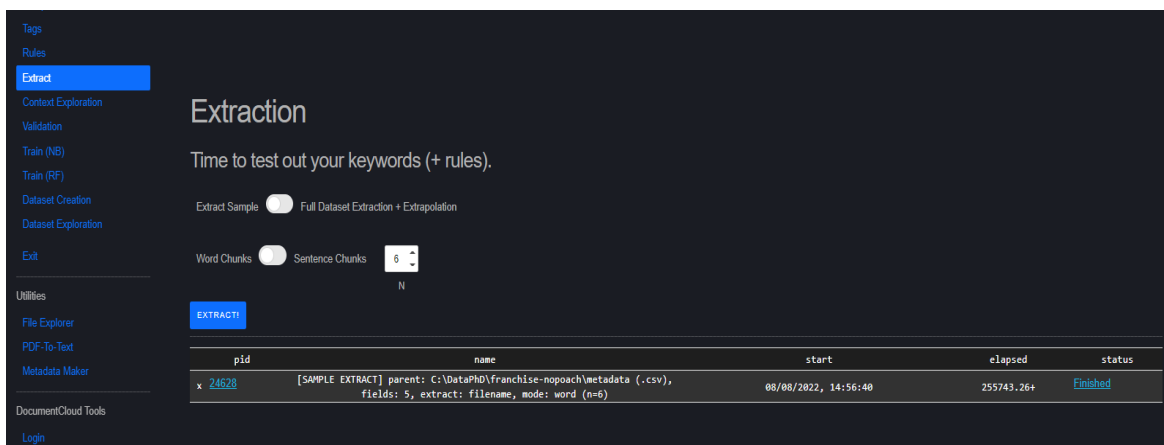
Keywords and tags are stored in keywords.json.

In CRAML, "Tags" appear on the left side of the screen. Keywords can be added in CRAML (click "Add Keyword" then "Save Keyword") or outside of the browser.
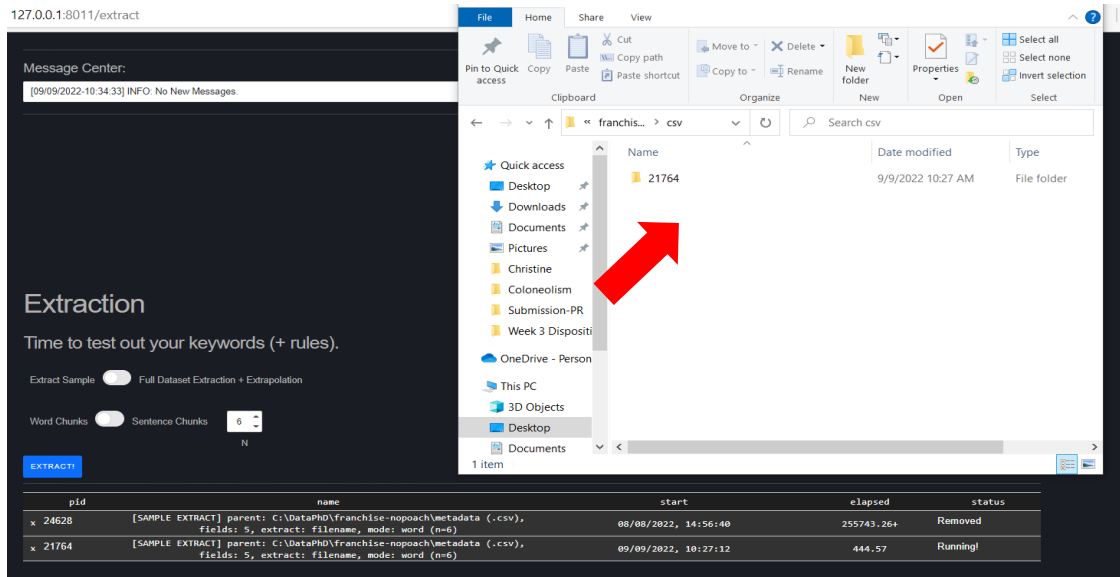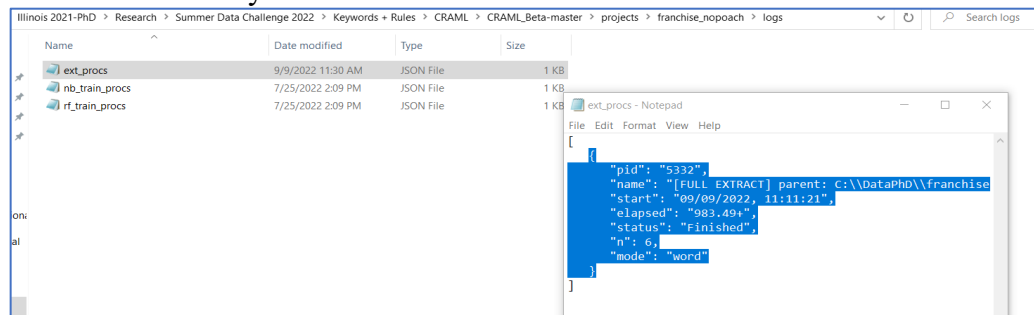
## Extracting chunks

a) Extract chunks of text containing keywords from the sample.

b) Choose "Word" or "Sentence" length chunks and adjust the N (number of words or sentences to the left and right of the keyword that will be extracted). As small an N as necessary to understand the context is advised.

c) Click "EXTRACT!" A folder with the process identification number for the extraction (pid) will appear in the projects folder under "csv".

d) The csv file in the pid folder will contain all extracted text from the sample.

e) Users who do not wish to create ML models can extrapolate the rules onto the full dataset by selecting "Full Dataset Extraction + Extrapolation."

f) The "logs" folder contains an "ext_procs.json" file that contains the information on each extraction. Manually erasing the pid as shown in the screenshot below will remove the entry in CRAML.



## Exploring context of extracted text data

a) "Context Exploration" allows users to identify the most common NGrams in the extracted text that contain keywords within a specific keyword.

## Rules

a) Rules are stored in csv files in the project/rules folder.

b) Each rule has a priority level and a 0/1 tag.

c) Keywords are given a priority level zero – meaning the first thing that CRAML is going to do is get all of the chunks of text that has the keyword and code it 0. Higher priority rules overwrite lower priority rules.

**Validation**

a) To understand the interaction between the rules and the extracted text, Validation allows uers to set a minimum parameter for how much observations per rule.

b) Users can validate within CRAML and/or export validation files to third party RAs to independently validate.



c) In CRAML, users can go over the chunks and see if they align with the rule



d) Validation files are saved in the projects under train/val

**Building Training Data and Training ML models in CRAML**

a) Under Train, select the sample, the validated rule file, and then click "train."
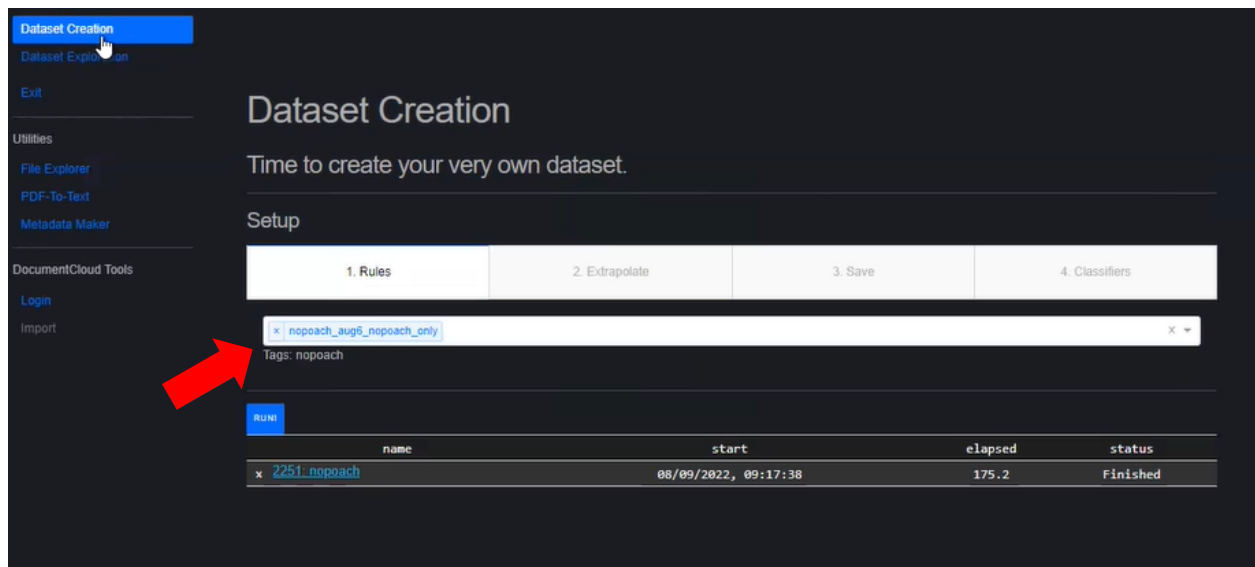


b) Users may need to augment resulting training data, either with several options in CRAML or outside. If the sample has a lot of "1s" but few "0s", of vice versa, it won't perform very well.

c) Clicking on "SHOW RESULTS" will report the F-Score and relevant information on the classifier produced.
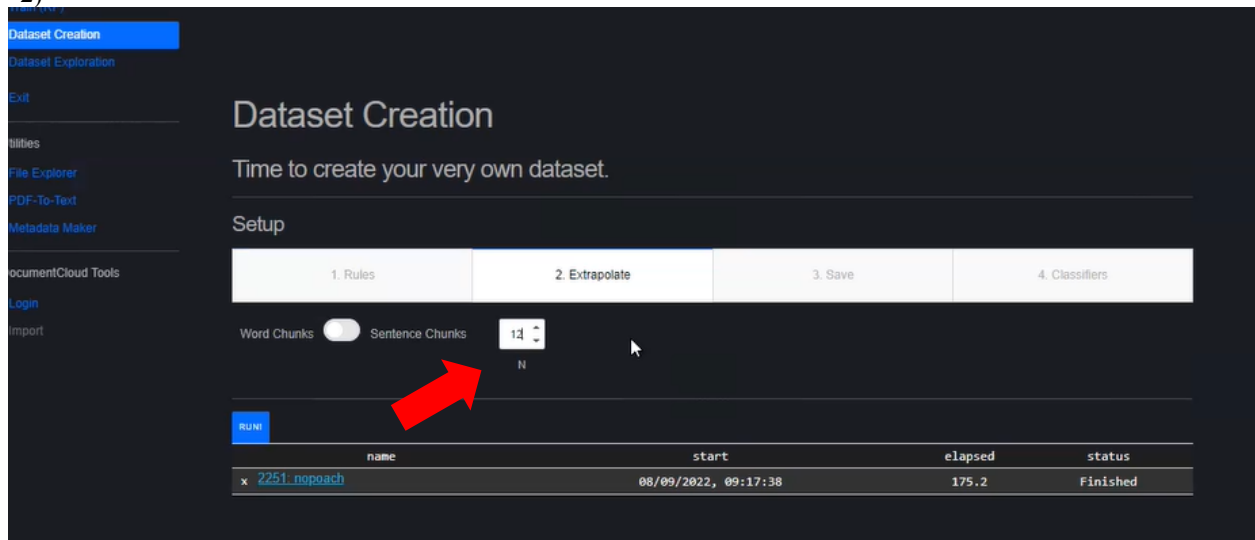
**Data Creation**

a) Creating a dataset from the machine learning model runs the entire pipeline of tools, and requires users:

a.  select a rules file,

b.  choose the context window,

c.  whether to keep the text or put it into a database,

d.  select a trained ML classifier,

e.  click "Run."



2)

3)