# Phase 1 Engine Documentation

The purpose of this documentation is to demonstrate the steps a user should take to successfully use the Engine, as well as interpret the outputs of the Engine. This documentation is for any potential user of the Engine. ==**The engine uses the Multilabel MLP**==.

## Instructions

1. Download and unzip project zip file: "*fyp_code_vX*", where *vX* is the version number.

   **1.1.** **Please ensure the folder has the following content before proceeding:**
   A. Engine
      - __pycache__
      - files
      - weights
      - app.py
      - predict.py
      - preprocess.py
      - sector_master_definition.xlsx
   B. .gitignore
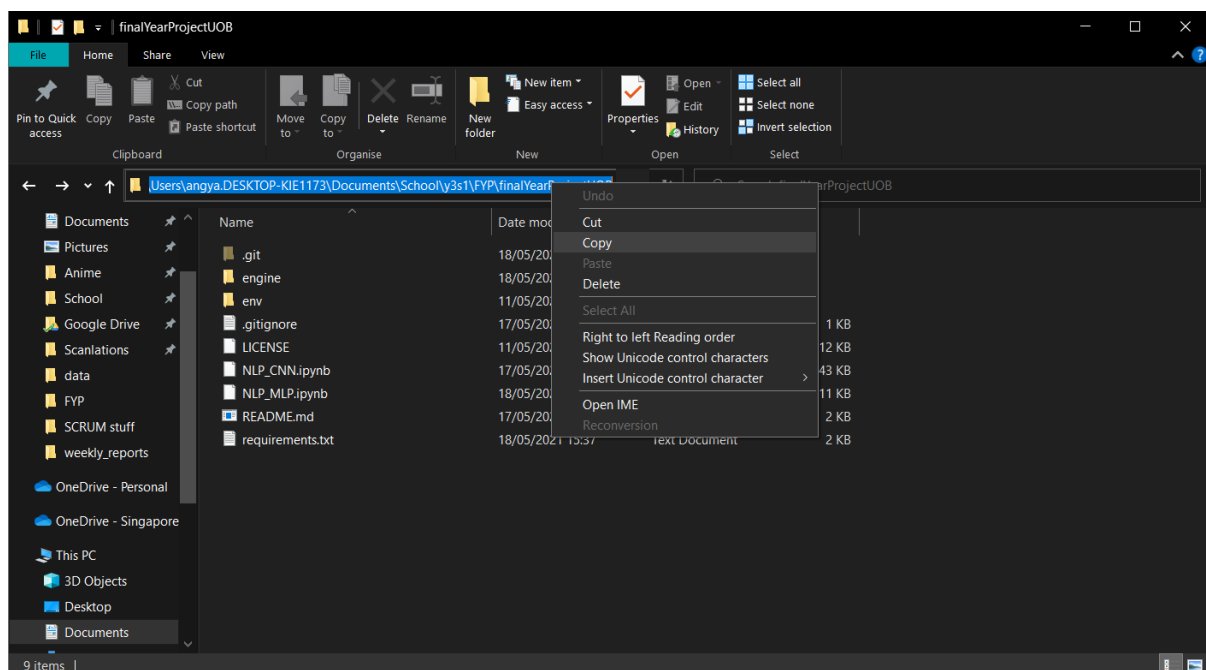   C. LICENSE
   D. NLP_CNN.ipynb
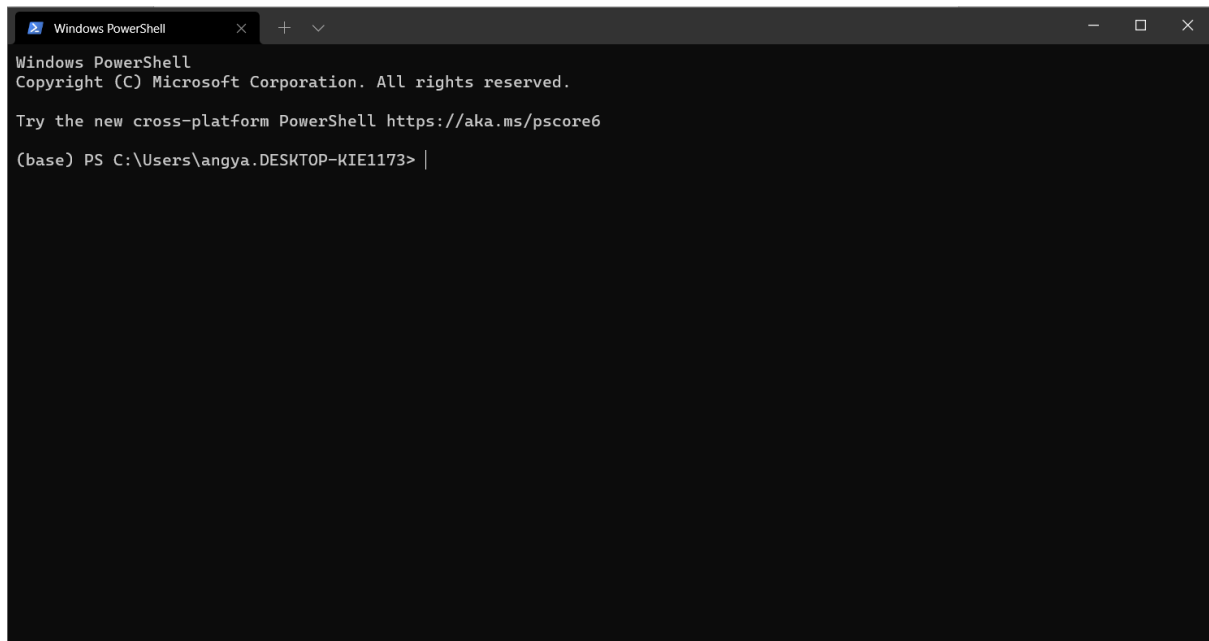   E. NLP_MLP.ipynb
   F. README.md
   G. requirements.txt

2. Navigate into unzipped project folder on your file explorer.
3. Copy path to project folder as such:

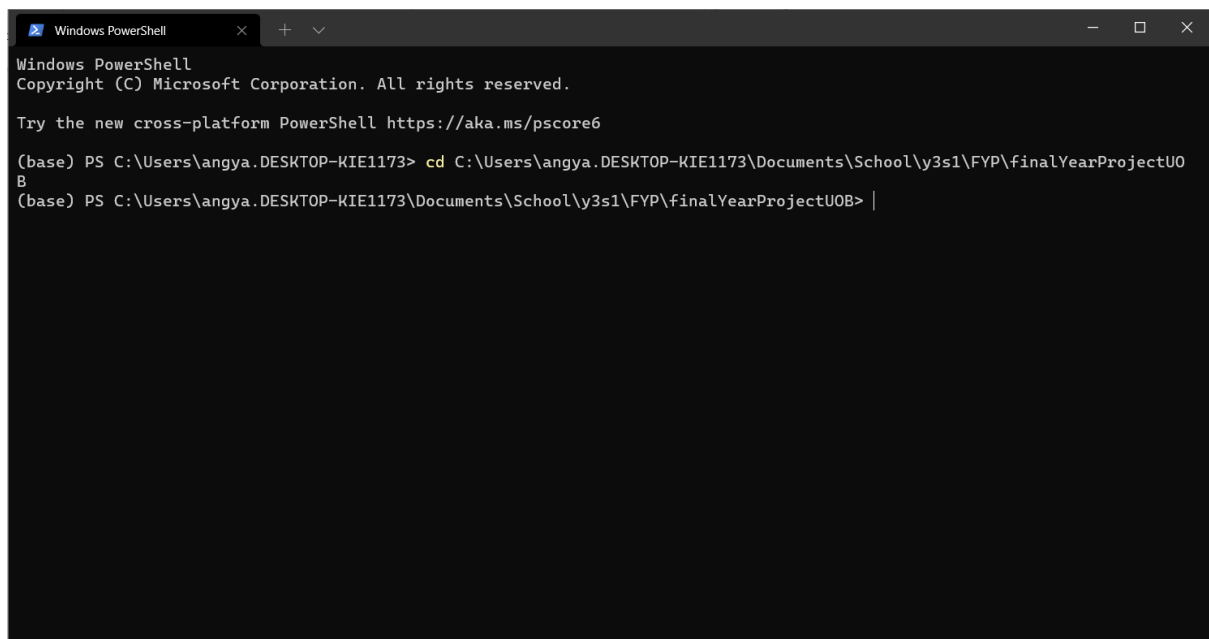4. Open Anaconda Prompt or Powershell.



5. Run `cd <project_path>`, where `<project_path>` is replaced with the copied path from the previous step.

6. Run `conda activate` to ensure that the Anaconda environment is enabled. (There should be a (base) visible if it is enabled)

7. OPTIONAL: Run `conda create -n temp python` to create a fresh conda environment, so that your default anaconda environment is not dirtied with the libraries required to run the code. Press y then Enter when prompted.

```
(base) PS C:\Users\angya.DESKTOP-KIE1173> conda create -n temp python
Collecting package metadata (current_repodata.json): done
Solving environment: done

## Package Plan ##

  environment location: C:\Users\angya.DESKTOP-KIE1173\anaconda3\envs\temp

  added / updated specs:
    - python


The following packages will be downloaded:

    package                    |            build
    ---------------------------|-----------------
    certifi-2020.12.5          |   py39haa95532_0         141 KB
    pip-21.0.1                 |   py39haa95532_0         1.8 MB
    python-3.9.4               |       h6244533_0        16.4 MB
    setuptools-52.0.0          |   py39haa95532_0         725 KB
    tzdata-2020f               |       h52ac0ba_0         113 KB
    wincertstore-0.2           |   py39h2bbff1b_0          15 KB
    ---------------------------------------------------------
                                           Total:        19.2 MB

The following NEW packages will be INSTALLED:

  ca-certificates    pkgs/main/win-64::ca-certificates-2021.4.13-haa95532_1
  certifi            pkgs/main/win-64::certifi-2020.12.5-py39haa95532_0
  openssl            pkgs/main/win-64::openssl-1.1.1k-h2bbff1b_0
  pip                pkgs/main/win-64::pip-21.0.1-py39haa95532_0
  python             pkgs/main/win-64::python-3.9.4-h6244533_0
  setuptools         pkgs/main/win-64::setuptools-52.0.0-py39haa95532_0
  sqlite             pkgs/main/win-64::sqlite-3.35.4-h2bbff1b_0
  tzdata             pkgs/main/noarch::tzdata-2020f-h52ac0ba_0
  vc                 pkgs/main/win-64::vc-14.2-h21ff451_1
  vs2015_runtime     pkgs/main/win-64::vs2015_runtime-14.27.29016-h5e58377_2
  wheel              pkgs/main/noarch::wheel-0.36.2-pyhd3eb1b0_0
  wincertstore       pkgs/main/win-64::wincertstore-0.2-py39h2bbff1b_0


Proceed ([y]/n)?
```

8. Wait until the setup is complete.
9. Run `conda activate temp` to activate the created anaconda environment. (There should be a (temp) visible if it is successfully enabled)



10. Run `conda install -c conda-forge cudnn` to install libraries required for GPU usage. Press y then Enter when prompted. Wait for installation to complete.

11. Run `python -m venv env` to make a new python virtual environment. This is to ensure that the installed python libraries do not dirty your python environment, and that there are no conflicts with existing packages that can cause the Engine to fail.



12. Run `.\env\Scripts\activate` to activate the created python virtual environment. (There should be a (env) visible.)

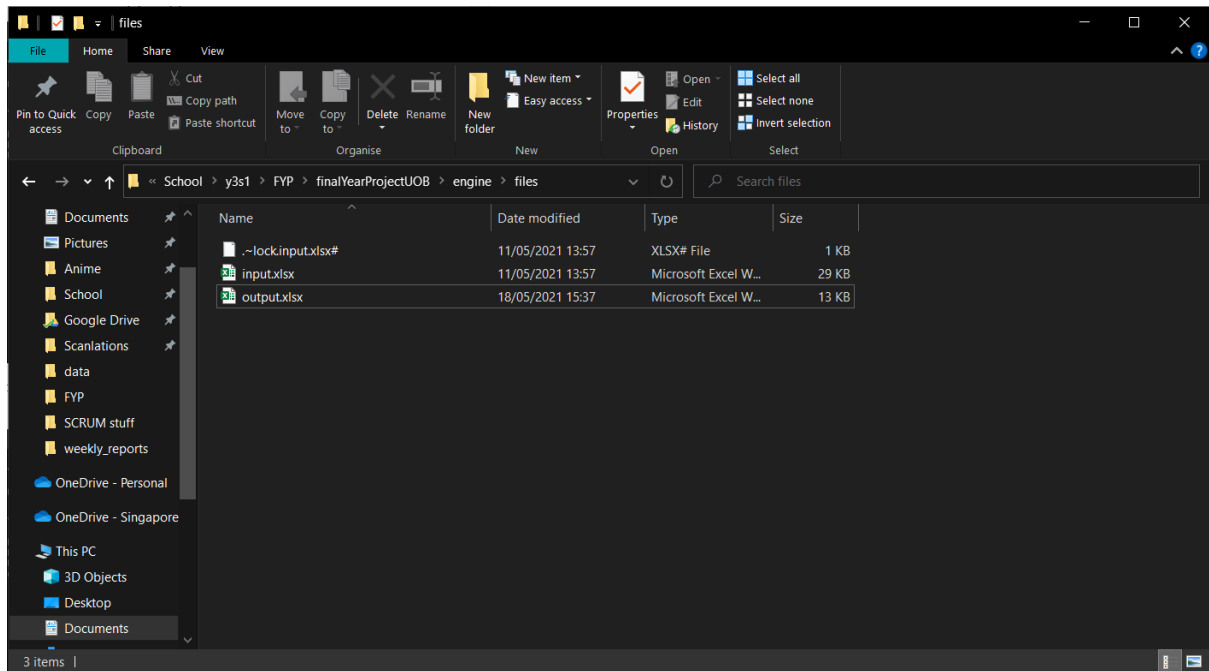13. Run `pip install -r requirements.txt` to install required libraries. Wait for installation to finish.

```
Collecting xlrd==2.0.1
  Using cached xlrd-2.0.1-py2.py3-none-any.whl (96 kB)
Collecting wheel<1.0,>=0.23.0
  Using cached wheel-0.36.2-py2.py3-none-any.whl (35 kB)
Requirement already satisfied: setuptools>=40.3.0 in c:\users\angya.desktop-kie1173\documents\school\y3s1\fyp\finalyearp
rojectuob\env\lib\site-packages (from google-auth==1.30.0->-r requirements.txt (line 17)) (49.2.1)
Installing collected packages: six, absl-py, wheel, astunparse, numpy, blis, cachetools, catalogue, certifi, chardet, cl
ick, fastrlock, cupy-cuda101, cymem, pydantic, murmurhash, pyparsing, packaging, tqdm, srsly, wasabi, typer, preshed, id
na, urllib3, requests, thinc, MarkupSafe, Jinja2, spacy-legacy, smart-open, pathy, spacy, en-core-web-lg, et-xmlfile, fi
lelock, flatbuffers, gast, pyasn1, rsa, pyasn1-modules, google-auth, oauthlib, requests-oauthlib, google-auth-oauthlib,
google-pasta, grpcio, h5py, joblib, scipy, PyYAML, Keras, keras-nightly, Keras-Preprocessing, Markdown, openpyxl, opt-ei
nsum, pytz, python-dateutil, pandas, protobuf, regex, sacremoses, spacy-alignments, spacy-lookups-data, typing-extension
s, torch, tokenizers, transformers, spacy-transformers, tensorboard-plugin-wit, Werkzeug, tensorboard-data-server, tenso
rboard, tensorflow-estimator, wrapt, termcolor, tensorflow, xlrd
Successfully installed Jinja2-2.11.3 Keras-2.4.3 Keras-Preprocessing-1.1.2 Markdown-3.3.4 MarkupSafe-1.1.1 PyYAML-5.4.1
Werkzeug-1.0.1 absl-py-0.12.0 astunparse-1.6.3 blis-0.7.4 cachetools-4.2.2 catalogue-2.0.4 certifi-2020.12.5 chardet-4.0
.0 click-7.1.2 cupy-cuda101-8.6.0 cymem-2.0.5 en-core-web-lg-3.0.0 et-xmlfile-1.1.0 fastrlock-0.6 filelock-3.0.12 flatbu
ffers-1.12 gast-0.4.0 google-auth-1.30.0 google-auth-oauthlib-0.4.4 google-pasta-0.2.0 grpcio-1.34.1 h5py-3.1.0 idna-2.1
0 joblib-1.0.1 keras-nightly-2.5.0.dev2021032900 murmurhash-1.0.5 numpy-1.19.5 oauthlib-3.1.0 openpyxl-3.0.7 opt-einsum-
3.3.0 packaging-20.9 pandas-1.2.4 pathy-0.5.2 preshed-3.0.5 protobuf-3.16.0 pyasn1-0.4.8 pyasn1-modules-0.2.8 pydantic-1
.7.3 pyparsing-2.4.7 python-dateutil-2.8.1 pytz-2021.1 regex-2021.4.4 requests-2.25.1 requests-oauthlib-1.3.0 rsa-4.7.2
sacremoses-0.0.45 scipy-1.6.3 six-1.15.0 smart-open-3.0.0 spacy-3.0.6 spacy-alignments-0.8.3 spacy-legacy-3.0.5 spacy-lo
okups-data-1.0.0 spacy-transformers-1.0.2 srsly-2.4.1 tensorboard-2.5.0 tensorboard-data-server-0.6.1 tensorboard-plugin
-wit-1.8.0 tensorflow-2.5.0rc3 tensorflow-estimator-2.5.0rc0 termcolor-1.1.0 thinc-8.0.3 tokenizers-0.10.2 torch-1.8.1 t
qdm-4.60.0 transformers-4.5.1 typer-0.3.2 typing-extensions-3.7.4.3 urllib3-1.26.4 wasabi-0.8.2 wheel-0.36.2 wrapt-1.12.
1 xlrd-2.0.1
WARNING: You are using pip version 20.2.3; however, version 21.1.1 is available.
You should consider upgrading via the 'c:\users\angya.desktop-kie1173\documents\school\y3s1\fyp\finalyearprojectuob\env\
scripts\python.exe -m pip install --upgrade pip' command.
(env) (temp) PS C:\Users\angya.DESKTOP-KIE1173\Documents\School\y3s1\FYP\finalYearProjectUOB>
```

14. OPTIONAL: If you encounter an OS error saying that the package could not be installed as there is no such file or directory when running the above command, it is most likely due to the file path being too long. To solve this, search up Registry Editor in your search bar and copy paste this path into the search bar of the Registry Editor found at the top: "Computer\HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\Control\FileSystem". Then, right click on "LongPathsEnabled" and change the "Value data" from 0 to 1. After this, you must close the Powershell or Anaconda Prompt, delete the env folder and restart from the very beginning.

** Please do note if you are worried that this step might affect your file system in the future, you can change it back to 0 after running the engine**

15. Ensure that there is an "input.xlsx" in the folder /engine/files in the project folder. An "output.xlsx" is optional.



16. OPTIONAL: If you want to create a custom input file, do follow format of the already existing "input.xlsx" file, and make sure that the custom file is in /engine/files. (If the file contains invalid entries, it will be automatically removed and will not show up in the "output.xlsx" file)

17. To run the Engine, simply run `cd engine`, then `python app.py` in Powershell or Anaconda Prompt. Below is a partial snapshot of the terminal output.
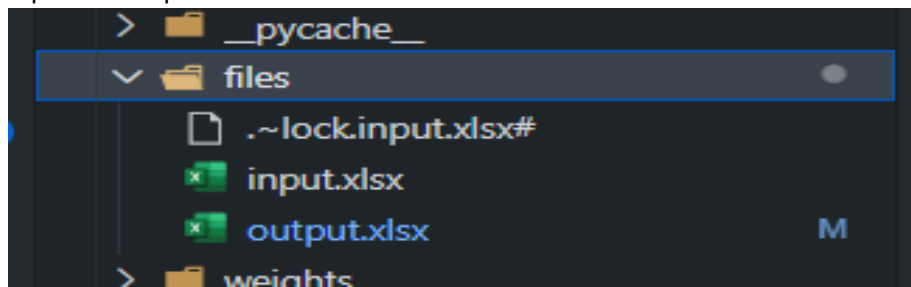




18. The output Excel file will be placed in the folder `/engine/files` and named "output.xlsx". Once opened, it will look like this:



19. To run the Engine with another Excel file, simply repeat steps 14 – 18.

## In-depth Descriptions

1. Input files: input.xlsx



      **input.xlsx** – The input file for the dataset for the engine to predict results.
      **Output.xlsx** – The output file from the engine, which includes the predicted results.
      The engine reads the input.xlsx file by feeding the data into app.py:

2. app.py



```python
import pandas as pd
from predict import predict

df_test = pd.read_excel('./files/input.xlsx')

# do prediction
df_test = predict(df_test)

# write results to excel file
df_test.to_excel('./files/output.xlsx', index=False)
```

Here, we read the dataset from the excel file and then populate it into a pandas dataframe. We then feed the data to the predict function found in predict.py. Finally when the prediction is done, we write the result to an output excel file.

1. predict.py

```python
1   # import necessary libraries
2   from tensorflow.keras.models import load_model
3   import pandas as pd
4   import numpy as np
5   from preprocess import preprocess
6
7   # load models
8   sector_model = load_model('engine/weights/model_1')
9   subsector_model = load_model('engine/weights/model_2')
10  archetype_model = load_model('engine/weights/model_3')
11  valuechain_model = load_model('engine/weights/model_4')
12
13  models = [sector_model, subsector_model, archetype_model, valuechain_model]
14
15  # import sector master definition file
16  df_keywords = pd.read_excel('engine/sector_master_definition.xlsx')
17
18  # preprocess dataset
19  df_keywords.drop(['Explanations', 'Notes'], axis=1, inplace=True)
20  df_keywords['Value Chain'].fillna(' ', inplace=True)
21  df_keywords['Sector Keywords'].fillna('[]', inplace=True)
22  df_keywords['Sector Keywords'] = df_keywords['Sector Keywords'].str.upper()
23  df_keywords.dropna(axis=0, how='any', inplace=True)
24
25  # read and split classification tags
26  sectors = sorted(list(df_keywords['Sector'].str.upper().unique()))
27  subsectors = sorted(list(df_keywords['Subsector'].unique()))
28  archetypes = sorted(list(df_keywords['Archetype'].unique()))
29  valuechains = sorted(list(df_keywords['Value Chain'].str.upper().unique()))
30
31  class_counts = [len(sectors), len(subsectors), len(archetypes), len(valuechai
32  classes = [sectors, subsectors, archetypes, valuechains]
33
34  # build keyword master list
35  keywords = []
36  for index, item in df_keywords['Sector Keywords'].iteritems():
37      keywords += eval(item)
38
39  keywords = sorted(list(set(keywords)))
40
41  # function for processing prediction results
42  def __process_results(result):
43      temp = []
44
45      for r in result:
46          temp.append((np.argmax(r), r[np.argmax(r)]))
47
48      return temp
```

```python
52  # === MAIN PREDICT FUNCTION === #
53  def predict(df):
54      df = preprocess(df, keywords)
55
56      X_pred = np.array(list(df['BoW_vectors']))
57
58      # do prediction
59      results = []
60      for model in models:
61          results.append(model.predict(X_pred))
62
63      # process output into rows
64      processed = []
65      for result in results:
66          processed.append(__process_results(result))
67
68      processed = np.array(processed)
69
70      # print results in human readable form
71      print('Prediction' + ' '*31 + '| Confidence')
72      print('-'*53 + '\n')
73
74      for index, row in enumerate(processed.swapaxes(0, 1)):
75          print('Company:', df['Company'].iloc[index], '\n')
76
77          for i in range(len(row)):
78              print(f'{classes[i][int(row[i][0])]}: <40.40} | {row[i][1] * 100: >9.4}%')
79
80          print('\n' + '-'*22 + '\n')
81
82      # add results to df
83
84      processed_tags = []
85      for i, result in enumerate(processed):
86          temp = []
87          for j, _ in result:
88              temp.append(classes[i][int(j)])
89
90          processed_tags.append(temp)
91
92      df['Sector'] = processed_tags[0]
93      df['Subsector'] = processed_tags[1]
94      df['Archetype'] = processed_tags[2]
95      df['Valuechain'] = processed_tags[3]
96
97      return df.drop(['BoW_vectors', 'processed'], axis=1)
```

At predict.py, what we did first was to load the pre-trained model weights that were done in google collab and also gather the sector master definition file which was already found in the master directory of the repository, and then build our master list keywords.

__process_results(): A function to preprocess the predicted results, getting the argmax, or the highest percentage recorded score.

predict(): A function whereby we can preprocess our raw data first by calling another function called preprocess() coded in another file. We then use the returned result from that function to format the result and print out the confidence score based on each company.

predict(): A function whereby we can preprocess our raw data first by calling another function called preprocess() coded in another file. We then use the returned result from that function to format the result and print out the confidence score based on each company.

2. preprocess.py

```python
# import libraries
import spacy
import pandas as pd
from spacy.language import Language
from spacy.tokens import Doc
from spacy.lang.char_classes import ALPHA, ALPHA_LOWER, ALPHA_UPPER, CONCAT_QUOTES, LIST_ELLIPSES, LIST_ICONS
from spacy.util import compile_infix_regex

# set up spacy
spacy.prefer_gpu()

# init nlp
nlp = spacy.load('en_core_web_lg')

# ----- passthrough ----- #
# keywords masterlist
#

# declare custom properties
Doc.set_extension('processed', default=True, force=True)
Doc.set_extension('word_bag', default=True, force=True)

# Modify tokenizer infix patterns
infixes = (
    LIST_ELLIPSES
    + LIST_ICONS
    + [
        r"(?<=[0-9])[+\-\*^](?=[0-9-])",
        r"(?<=[{al}{q}])\.(?=[{au}{q}])".format(
            al=ALPHA_LOWER, au=ALPHA_UPPER, q=CONCAT_QUOTES
        ),
        r"(?<=[{a}]),(?=[{a}])".format(a=ALPHA),
        r"(?<=[{a}0-9])[:<>=/](?=[{a}])".format(a=ALPHA),
    ]
)

infix_re = compile_infix_regex(infixes)
nlp.tokenizer.infix_finditer = infix_re.finditer
            dictionary = dict.fromkeys(keywords, 0)
            for word in company:
                if word in keywords:
                    dictionary[word] += 1

            # append to dataframe
            bow_vectors.append(list(dictionary.values()))

            # print(f'{sum(dictionary.values()):>3}/{len(dictionary.values()):<3} |', dictionary.values())

        df['BoW_vectors'] = bow_vectors

        return df
```

At preprocess.py, where the data will be fed in from the predict.py after getting called, we will feed the data into spacy to do preprocessing, which includes tokenization, lemmenization, word cleaning, removing stop words, and punctuations.

## 3. Results

Once the entire process has finished running, users can see the following results:

```
Prediction                       | Confidence
--------------------------------------------------

CNI                              |    99.99%
cni_service providers            |    99.96%
cni_service providers            |    99.18%
MIDSTREAM                        |    100.0%


---------------------

CNI                              |    69.4%
buildings & industrial           |    52.23%
buildings & industrial_contractor|    44.03%
MIDSTREAM                        |    64.41%


---------------------

CNI                              |    100.0%
building material                |    100.0%
building material_manufacturer   |    100.0%
MIDSTREAM                        |    100.0%


---------------------

CNI                              |    100.0%
buildings & industrial           |    100.0%
buildings & industrial_contractor|    100.0%
MIDSTREAM                        |    100.0%


---------------------
```

And the output in the excel file: