This project is developed based on Python 3.9. Therefore, when downloading Python, please ensure its version is >=3.9. You can download Python from this link: https://www.python.org/downloads/ We highly recommend using a Python IDE to run our project. We suggest using PyCharm. You can download it from this link: https://www.jetbrains.com/pycharm/
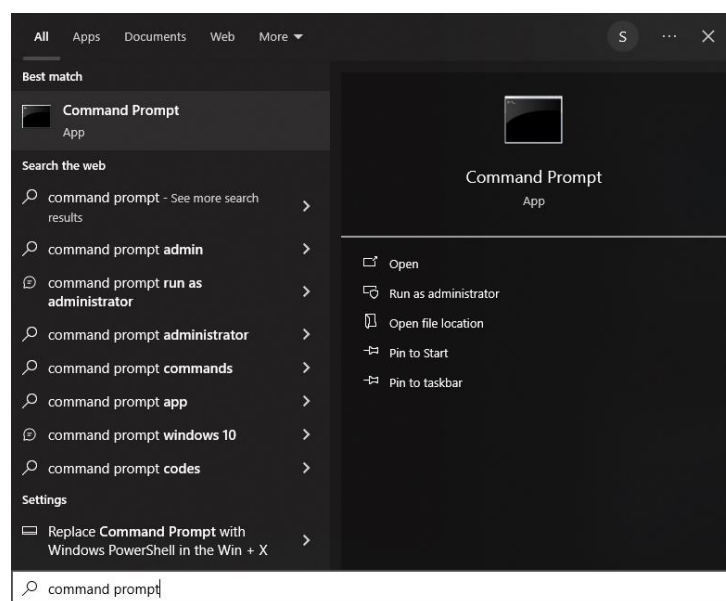1;

## How to download all 300 FELO Leaflets (Webscraping)：

*Tested on Windows 10

### Pre-requisites:

1.  Download Google Chrome
2.  Have Python 3 installed
3.  Have the necessary packages installed. To install the packages, open Command Prompt via Start. without quotation marks, run '`pip install requests beautifulsoup4 selenium`'
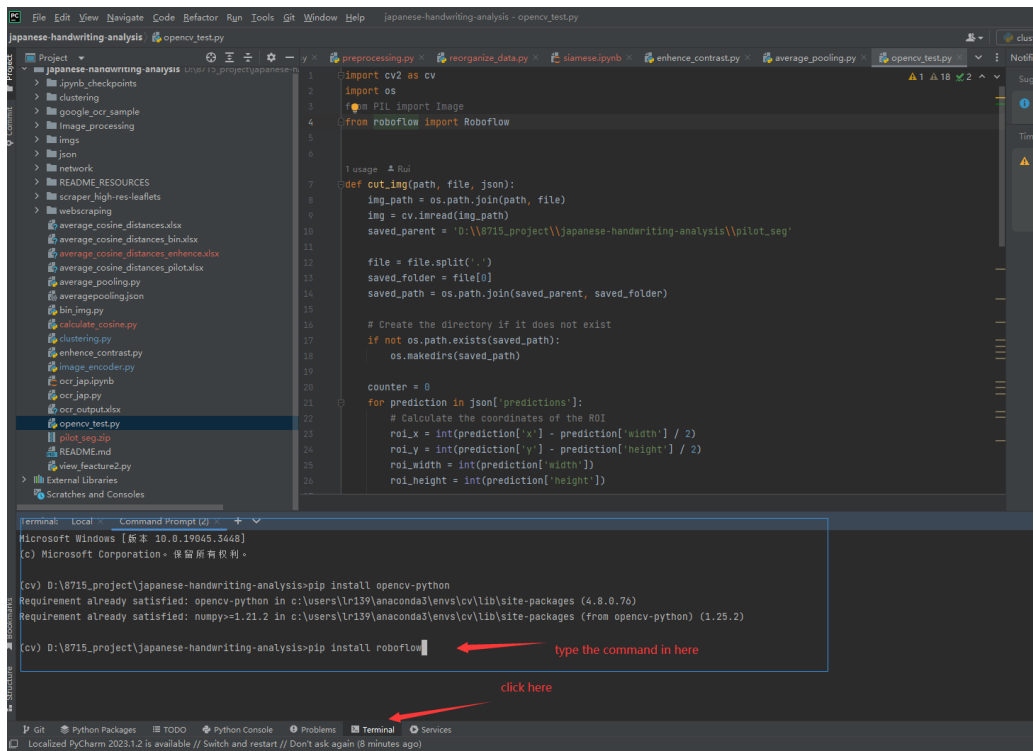


### How to download the leaflets:

1.  Within the Command Prompt, navigate to /scraper_high-res-leaflets/ folder.
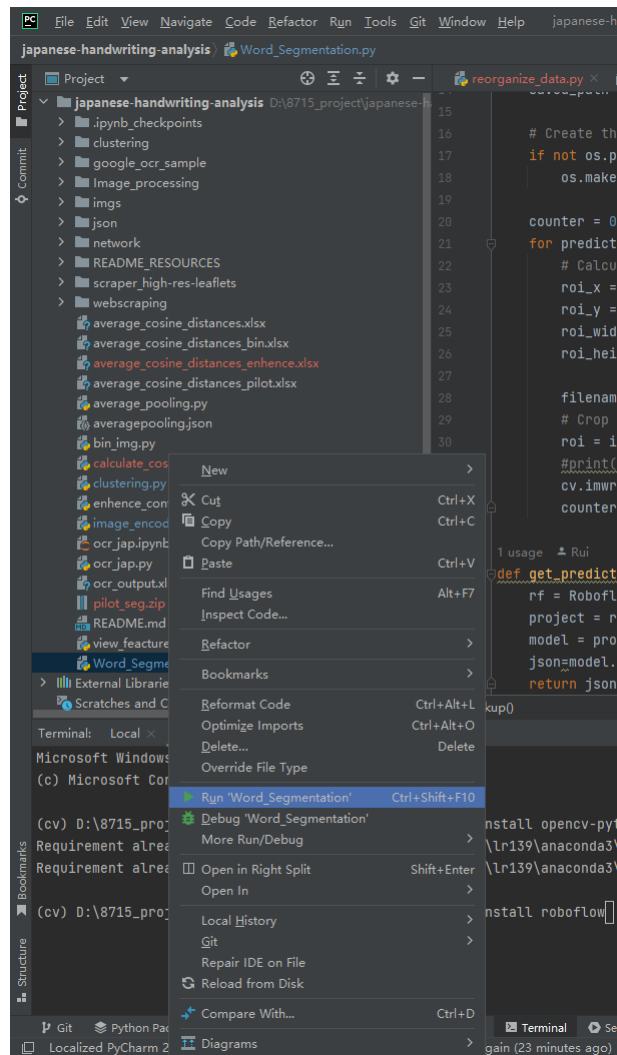2.  Without quotation marks, run '`python scraper.py `' or '`python3 scraper.py`'

Leaflets will be downloaded within the same directory the script is ran. Based on FELO_leaflets spreadsheet, there should be 300 leaflets downloaded. Download speed may vary depending on internet speed.

# Word Segmentation：

- Use the following commands to install the required libraries
  pip install opencv-python
  pip install roboflow



- Open the Word Segmentation.py file. In this file, we will invoke the model from RoboFlow to implement Word Segmentation. What you need to modify are:

  - Line 10's saved_parent: This is the path where the output images will be saved.
  - Line 42's path: This is the path for the input images.
  - Line 36's api key: Refer to the RoboFlow section for this

- The path for the input images should be in the following format:
  - -japanese-handwriting-analysis
  - ----input_image_folder
  - ----------img1.jpg
  - ----------img2.jpg
  - ----------….jpg
  - ----other_folder
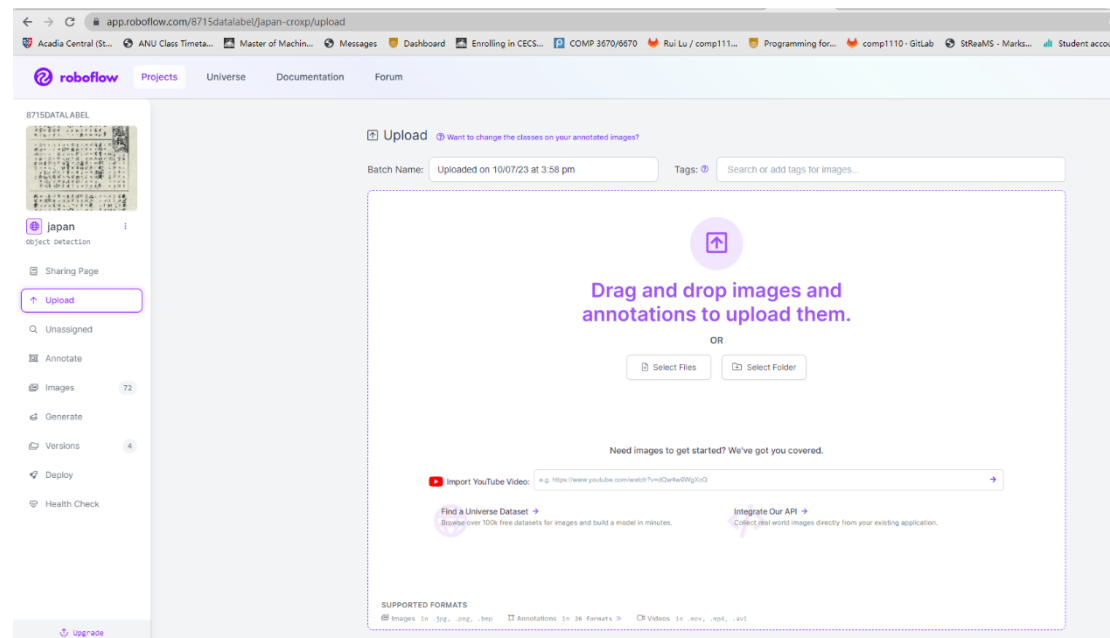  - ----Word_Segmentation.py

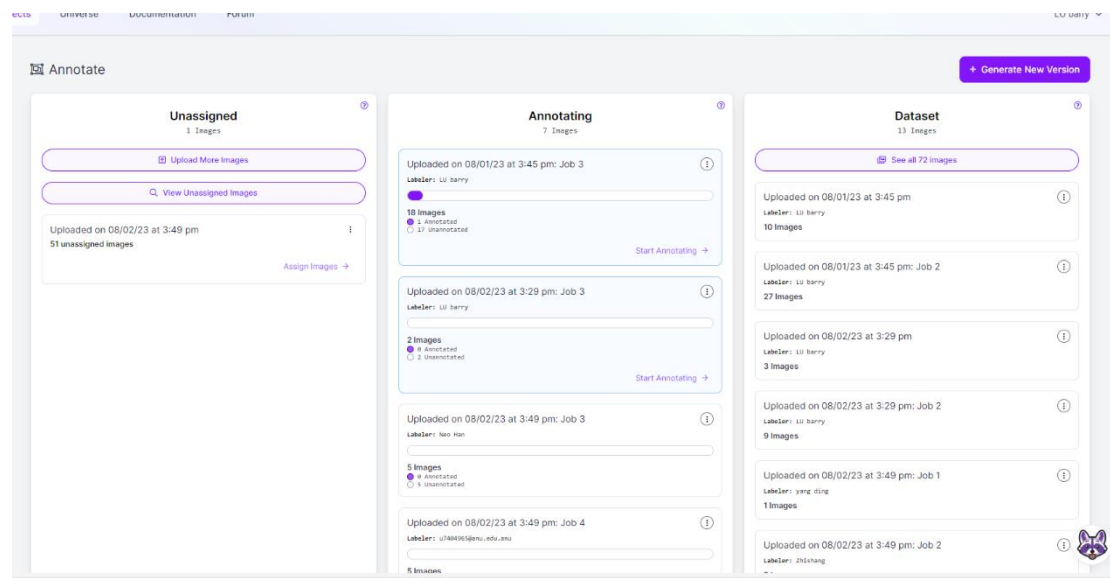- In the left navigation pane, select the file and right-click, then choose Run

## 2: RoboFlow

RoboFlow offers an integrated platform for annotating data, training models, and deploying models. Our current RoboFlow project is nearing its API call limit. Therefore, you might consider creating a new RoboFlow project and transferring the annotated data into it. The annotated data can be found on Google Drive under the filename 'japan.v5i.yolov8'.

Train model using RobowFlow:



First, open your RoboFlow project, then go to the 'upload' page to upload your dataset



If you want to annotate more data, go to the 'annoted' page. Click on 'assign image' to allocate the images to yourself or your team members. Then, click 'start annotating'. After completing the annotations, click on 'generate new version'.

Follow the instructions on this page, and then you can train your model.



After completing the training, click 'deploy', and you can then invoke your model. You can modify the content in Word Segmentation.py based on the Python code they provide.

## 3: Clustering


## OCR branch:

See the full plan in research report on confluence [Report on OCR method progress, plan and constraints - International Collaboration on Handwriting - Confluence (atlassian.net)](#)
Still lacking training and testing. Only finished finding OCR and use the model.
The result of seg_letter is in ocr_output.xlsx, if you want to use the model again. You can go to clustering/google_ocr.ipynb. Run every cell except last one, remember to change the input and output directory! The seg_letter is on google drive so probably download it to your computer and lead the path to it.
Ask Zhishang Bian if you have any questions.


## Cluster with K-means:


This is work from semester 1 2023. New methods don't use it now. You probably don't need it
Open with Pycharm or Jupyter lab.
Jupyter lab: Run every cell and it will automatically download the library you need
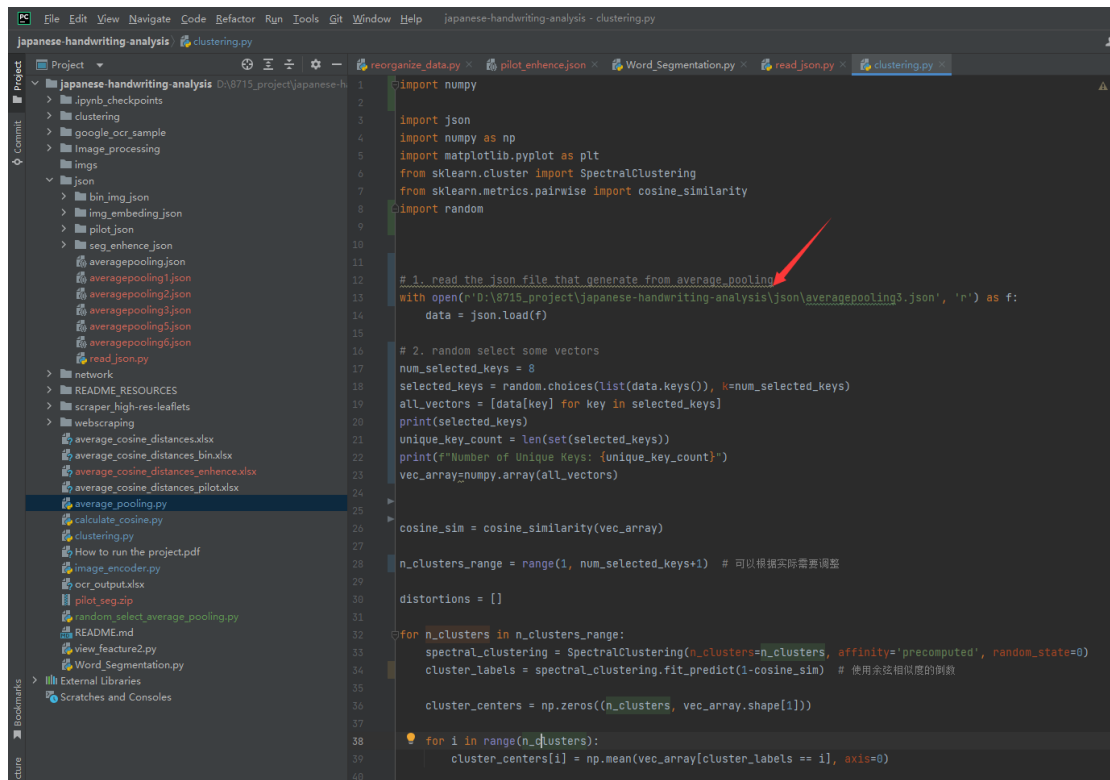Pycharm: Run main() and it will automatically download the library you need
Found it here:
[https://gitlab.cecs.anu.edu.au/u7434576/japanese-handwriting-analysis/-/tree/main/clustering](https://gitlab.cecs.anu.edu.au/u7434576/japanese-handwriting-analysis/-/tree/main/clustering)


# Spectral clustering：

First,run the Word Segmentation.py,image_encoder.py,averagepooling.py

Modify the path on line 13. You can generate the json file yourself or use the one in gitlab. Then you can run the program

## 4：Image encoder：

Install the following package pytorch, pil,clip

The segamentaed image can be downloaded from here:

https://drive.google.com/drive/folders/1hecluxT4ad-6tQBP7q64SHWG3DRKsFKW

seg_letter.zip is the leaflets. Pilot_seg,zip is the test data. You can ou can unzip them in the root directory.

```python
def get_image_features(model, preprocess, folder_path):
    image_features_dict = {}

    for subdir, _, files in os.walk(folder_path):
        for file in files:
            if file.endswith(('.png', '.jpg', '.jpeg')):
                # get_path
                full_path = os.path.join(subdir, file)

                # use clip to get the feactures
                image_tensor = preprocess(Image.open(full_path)).unsqueeze(0).to(device)
                with torch.no_grad():
                    image_features = model.encode_image(image_tensor)
                image_features = image_features.cpu().numpy().tolist()

                # save feacture vector
                image_features_dict[full_path] = image_features[0]

    return image_features_dict


if __name__ == "__main__":
    device = "cuda" if torch.cuda.is_available() else "cpu"
    model, preprocess = clip.load("ViT-B/32", device=device)

    folders = {

        r"D:\8715_project\japanese-handwriting-analysis\seg_letter_enhance": "D:\8715_project\japanese-handwriting-ana
    }

    for input_folder, output_folder in folders.items():
        image_features_dict = get_image_features(model, preprocess, input_folder)
        # save to json format
        if not os.path.exists(output_folder):
            os.makedirs(output_folder)
        output_file = os.path.join(output_folder, os.path.basename(input_folder) + ".json")
        with open(output_file, 'w') as f:
            json.dump(image_features_dict, f, indent=4)
```

*image path* ↘

*output path* ↘

Modify the path on line 33. Then you can run it.

## 5: AveragePooling

```python
1 usage    Rui *
def process_json(input_path, output_path):
    with open(input_path, 'r') as f:
        data = json.load(f)

    # 结果数据初始化
    result_data = {}

    # key is the path of the img, values are the vectors
    dir_vectors = {}
    for img_path, vector in data.items():

        dir_name = os.path.basename(os.path.dirname(img_path))

        if dir_name not in dir_vectors:
            dir_vectors[dir_name] = []

        dir_vectors[dir_name].append(vector)

    for dir_name, vectors in dir_vectors.items():
        # select 30 vectors
        selected_vectors = vectors[:30]


        avg_vector = average_pooling(selected_vectors)


        result_data[dir_name] = avg_vector

    # save to json format
    with open(output_path, 'w') as f:
        json.dump(result_data, f)


input_path = 'D:\8715_project\japanese-handwriting-analysis\json\pilot_json\pilot_enhence.json'
output_path = 'json/averagepooling3.json'
process_json(input_path, output_path)
```

Modify the path on line 45 and 46. Then run the file.