# Group 3

Ebram Thabet - Omar Elwaliely - Mohamed Alashkar - Andrew Aziz

Documentation for web scraping the AUC Digital library
"https://digitalcollections.aucegypt.edu/digital/collection/p15795coll20/search"

Firstly, the user should run the Python script named "Webscrap_Selenium". This script will open the library website: "https://digitalcollections.aucegypt.edu/digital/collection/p15795coll20/search". It is very important **not** to close this website while the script is running.

For each book on this website, a new folder will be created in the directory where the Python script is being executed. The folder name will match the book's name. All images **containing faces** from the book will be downloaded into this folder as "jpg" files. The image filenames will follow the format: "BookName_PageName".

If two images have the same name, the first image will be saved as "BookName_PageName", and the second will be saved as "BookName_PageName (1)", with subsequent images following this numbering pattern.

If a folder is named "Unknown_Book_Name" or an image is named "Unknown_Book_Name_Unknown_Page_Name", it means the code could not find the book name or page name. The user may need to rerun the code to resolve this issue.

Books will be processed and downloaded in the order they appear on the digital library website. If the user stops the script and reruns it, the script will start from the beginning, iterating through the books in order. However, if it finds a folder with the same name as a book, it will skip that book and continue to the next one. This process will repeat until it reaches a book for which no corresponding folder exists. Essentially, the script resumes the downloading process from the last book it was working on before being stopped.

This code relies on the **HTML** structure of the AUC Digital Library website. If there are any changes to the website's HTML, the code may stop working or crash.

This Google Drive folder contains the **expected output** that the user should receive after running the whole code (Use AUC email to be able to access it): "https://drive.google.com/drive/folders/1NAi29f37pdWhIwIJSzQ6YNC1V4s1YS56?usp=sharing"

For the **OCR** (Optical Character Recognition) functionality, we use the library "Pytesseract" to extract text from images in each book. The OCR process specifically targets names written in capital letters or those preceded by titles such as 'Mr.', 'Mrs.', or 'Dr.'.

To use the OCR code, the user needs to specify the path to the folder containing the images by setting a variable named "folder_path". Once the path is set, the OCR code can be executed to perform the extraction and matching.
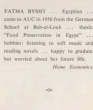
When the user runs the OCR code, it matches the extracted names with the correct face encodings in some images (you can find these images below) from the first book in the AUC digital library, titled 'Al Hawdaj of the Campus Caravan 1958'. This process correctly associates names with their corresponding face encodings and saves this data in a CSV file named 'faces_data.csv'.

Our GitHub repository also includes a CSV file named "Webscraped_Book.csv", which contains the names and face encodings for the below images from the book 'Al Hawdaj of the Campus Caravan 1958'. Below is a list of images that the current OCR can process and match with the corresponding face encodings:

ESTREGO NASSAR ... Palestinian ... *Humanities-History.*

EVA EL ZIK ... *Education.*

FARIDA TARAZI ... Palestinian ... *Home Economics.*

FAROUK ABOUL RAHIM ... Jordanian ... *Chemistry.*

FAROUK OUDDOUMI ... Palestinian ... *Economics.*

FATMA BYBIT ... Egyptian ... *Home Economics.*

FAUSTINO BORO ... *Journalism.*

FUAD BAKHIT ... Saudi Arabian ... *Economics-Business.*

GHADA EL RAFEY ... Lebanese ... *Political Science.*

GHALIB HALASA ... Jordanian ... *Journalism.*

GILBERT DOSS ... Egyptian ... *Economics-Business.*

HAMDEYA HAMDY ... Egyptian ... *Literature.*

LAURICE SIDHOOM ... *Sociology.*

LOIS GIRGIS ... Egyptian ... *Education.*

MAHMOUD EL TUKHI ... Egyptian ... *Education.*

MAHA ABU EL MONEM ... Egyptian ... *Political Science.*

Mrs. MARIE GIRGIS ... Egyptian ... *Literature.*

MARO BOGHIKIAN ... Egyptian ... *Psychology and Education.*

NADIA NAGUIB ... joined AUC ... Egyptian ... *Education.*

MAGDA TADROS ... Egyptian ... *Journalism.*

MAHMOUD BAHLOUL ... Bahraini ... *Political Science.*

MADIHA ABUL FUTUH ... Egyptian ... *Sociology.*

MIREILLE MORDO ... Greek ... *Economics-Business.*

MOHAMMAD MEDHAT HUSSEIN ... Egyptian ... *Economics-Business.*

NADIA EL SHAZLY ... Egyptian ... *Journalism.*

NADIA ARAFA ... Egyptian ... *Literature.*

NAJWA IMAM ... Palestinian ... *Literature.*

NASSER EL UBI ... Palestinian ... *Economics-Business.*

NAWAL AL RIFAI ... Egyptian ... *Physics-Mathematics.*

NORA EL KADY ... Egyptian ... *Literature.*

PANDELIN HALAMANDARIS ... Greek ... *English Literature.*

RAGHER JADOUNI ... Palestinian ... *Journalism.*

RASHAD KAMAL ... came to AUC ... *Psychology.*

RAYMOND TADROS ... Palestinian ... *Economics-Business.*

REGINA PRIEDHO ... Egyptian ... *Psychology.*

RUQAYYA ABDEL HALIM ... Egyptian ... *Home Economics.*

SABA ALARJA ... Jordanian ... *Chemistry.*

SUHAIL SABANEGH ... *Economics-Business.*

SALIMA SIGARI ... Jordanian ... *Sociology.*

SALWA MOHARIB ... Egyptian ... *Literature.*

SAMI KHAWAJA ... *Chemistry.*

SAMI HANNA ... came to AUC ... *Economics-Business.*

SAMIA HANNA ... Egyptian ... *Physics-Mathematics.*

SAMHA BIBAWY ... Egyptian ... *Home Economics.*

SIHAM MITRI ... Lebanese ... *Physics-Mathematics.*

SAMIRA MOSAD ... Egyptian ... *Journalism.*

SAMIRA GHORBIAL ... Egyptian ... *Literature.*

SAMIRA KIHOLLOS ... Egyptian ... *Literature.*