# User Manual

This user manual provides some suggestions on programs to install before trying to run the programs in this repository. Command line instructions for running each program are found at the top of each program in the repository.

# 1 Image-Processing Pipeline and Ensemble Programs

To run the ensemble and image-processing programs, you will need to install

1. Python3 (`https://www.python.org/download/releases/3.0/`)

2. Tesseract (`https://github.com/tesseract-ocr/tesseract/wiki`)

3. pytesseract (`https://pypi.org/project/pytesseract/`)

4. Google Cloud Client Library (`https://cloud.google.com/vision/docs/libraries#installing_the_client_library`)

5. Python Imaging Library (`https://python-pillow.org/`)

6. NumPy (`https://www.numpy.org/`)

7. glob (`https://docs.python.org/3/library/glob.html`)

8. os (`https://docs.python.org/3/library/os.html`)

9. argparse (`https://docs.python.org/3/library/argparse.html`)

The eight traineddata files should be placed in your 'tessdata' directory. The exact directory will depend both on the type of training data, and your

Linux distribution. Possibilities are `/usr/share/tesseract-ocr/tessdata` or `/usr/share/tessdata` or `/usr/share/tesseract-ocr/4.00/tessdata`.

You will also need to set up a Cloud Vision API project in the Google Cloud Platform Console. (`https://cloud.google.com/vision/docs/quickstart#set_up_a_google_cloud_vision_api_project`). In particular, you should replace line 13 in the Google/Tesseract ensemble program with your own Google Credentials.

## 2   Testing

To obtain the character accuracy for the ensemble output, you should fork this GitHub repository `https://github.com/Shreeshrii/ocr-evaluation-tools`.