

Smart OCR for Document Digitization

Introduction:

Overview:

1. We provide input through post method in web page.
2. Preprocessing of pdf occurs, within pytesseract. Text localisation, detection, segmentation, recognition.
3. It returns content in text format along with the confidence, bounding boxes for what it has been recognized.
4. Take that output and saving it in a file.

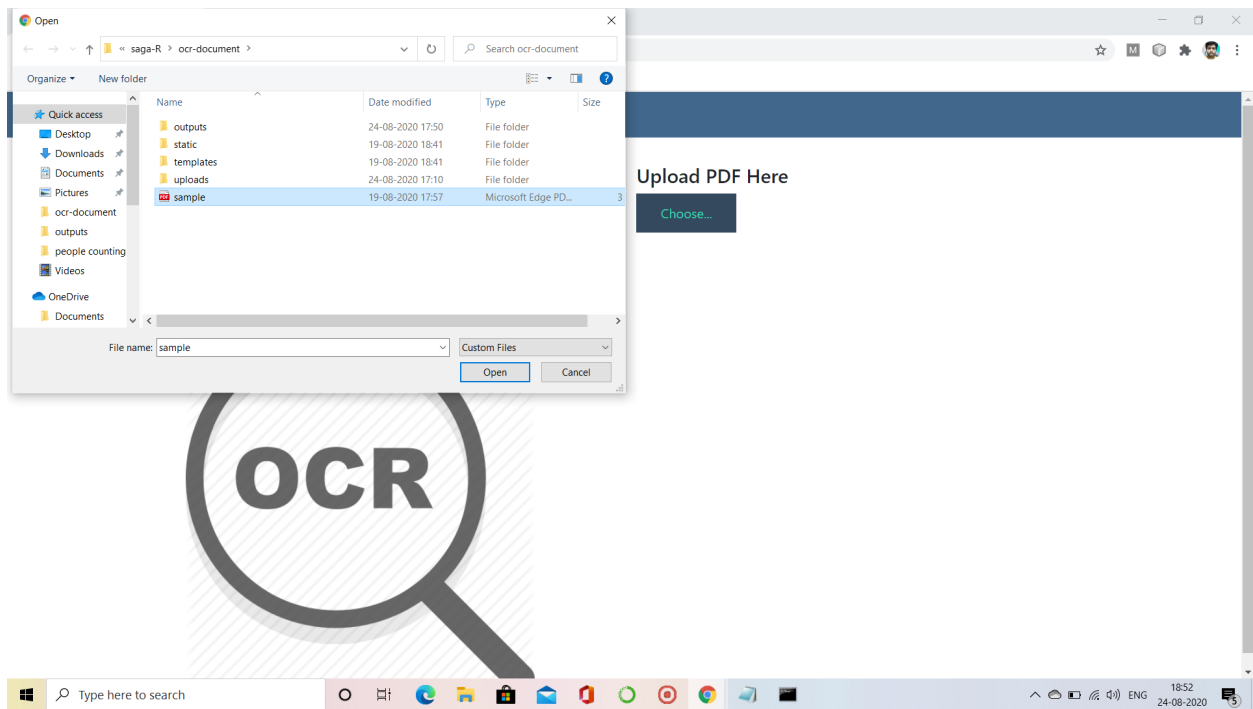
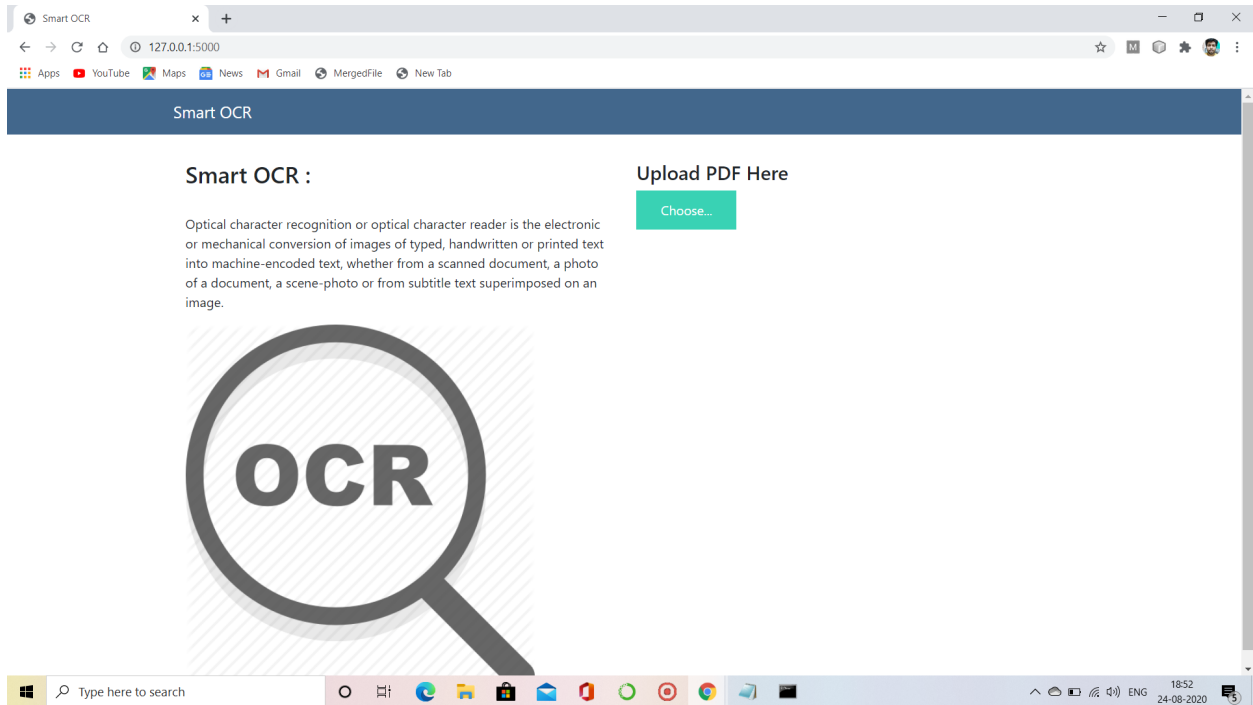
The whole process is dumped in Flask.

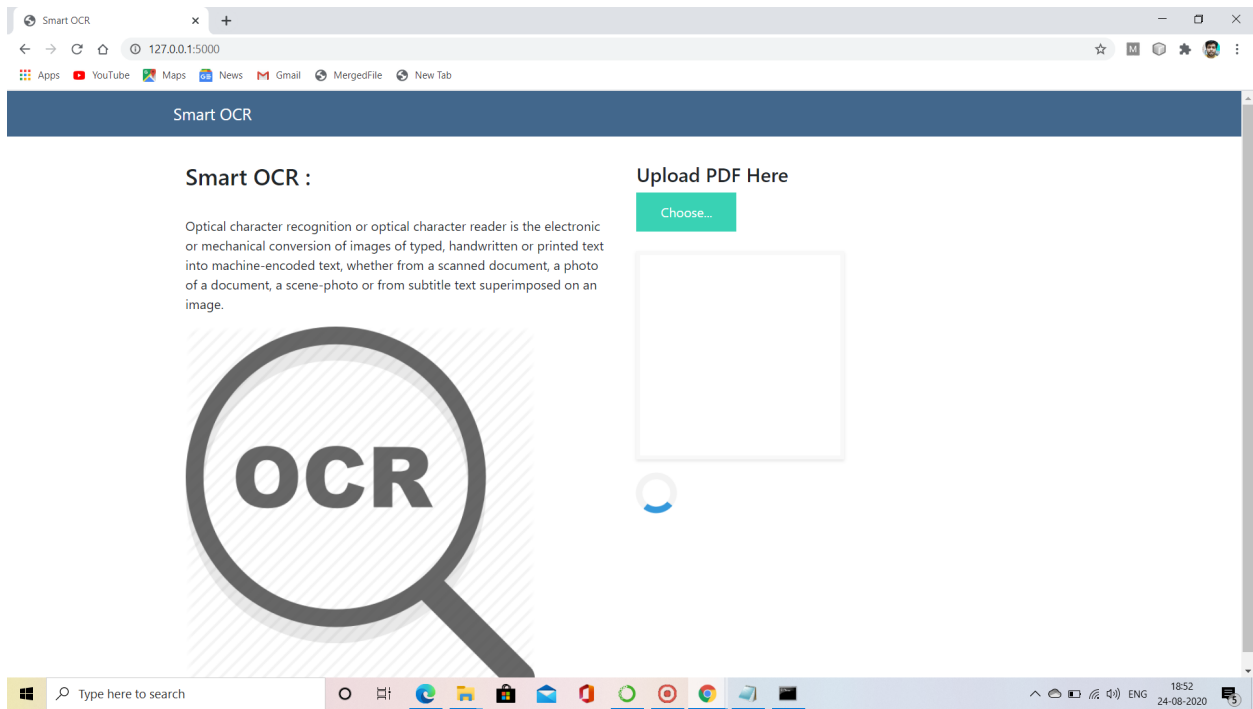
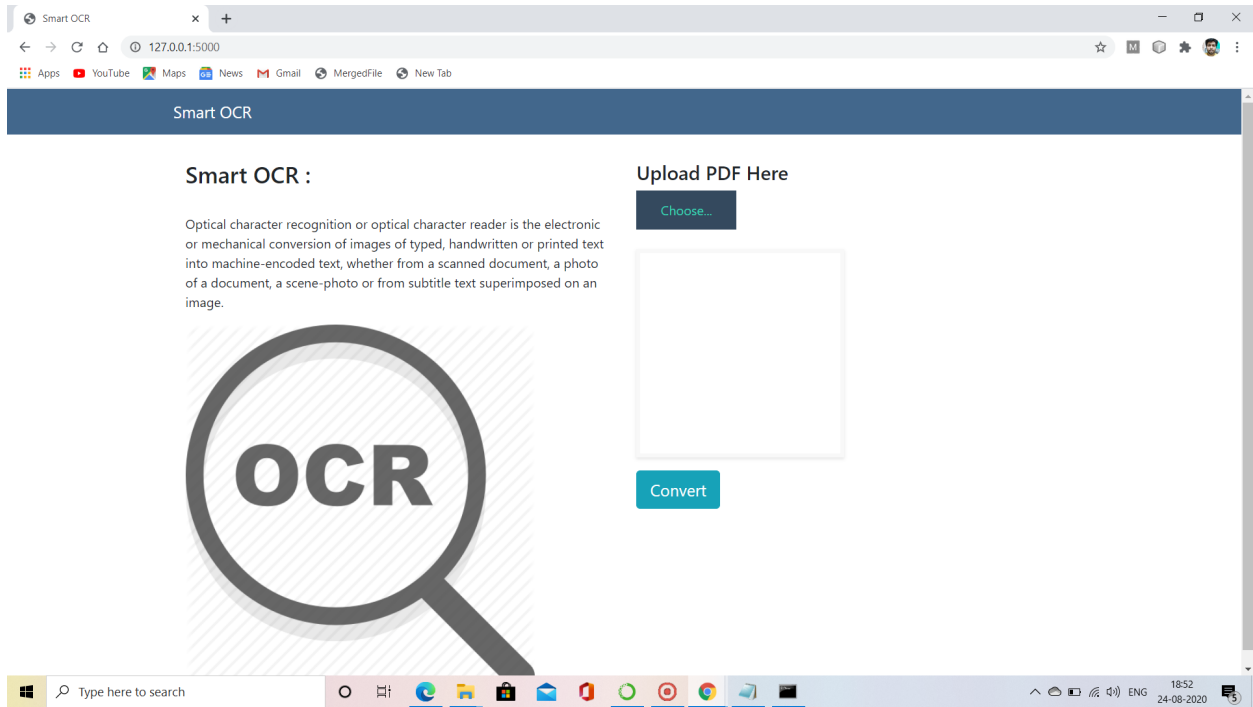
Purpose:

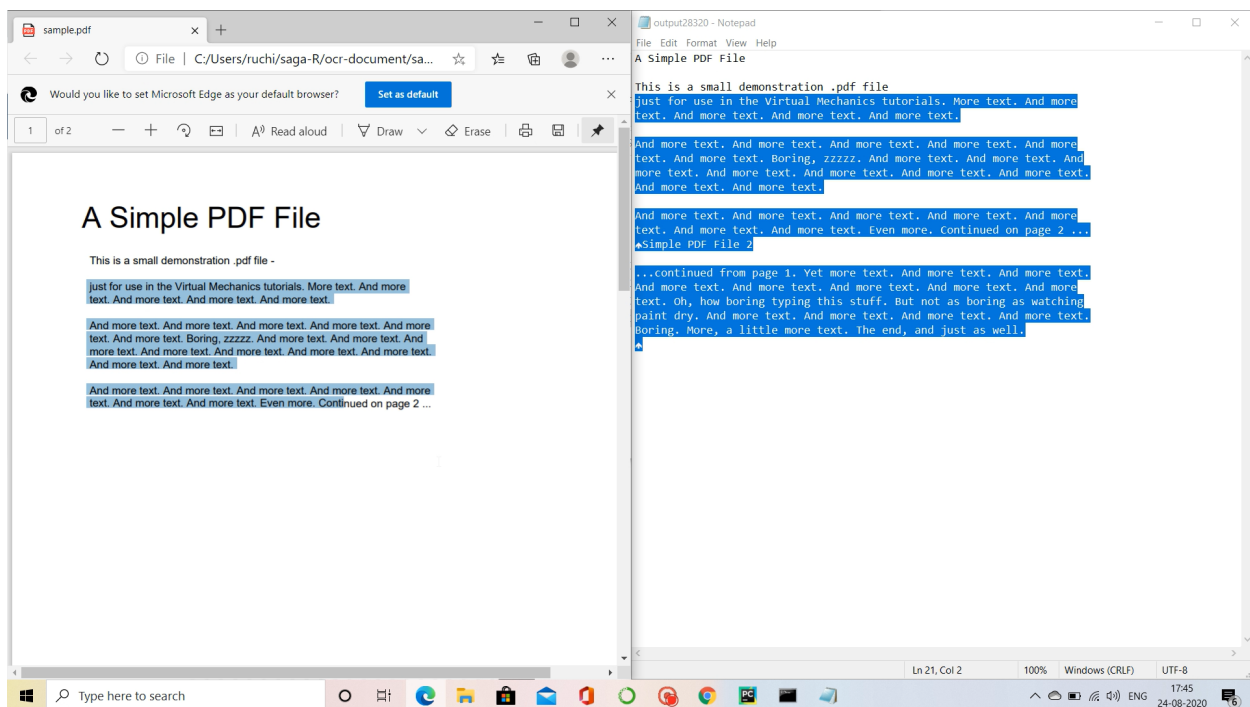
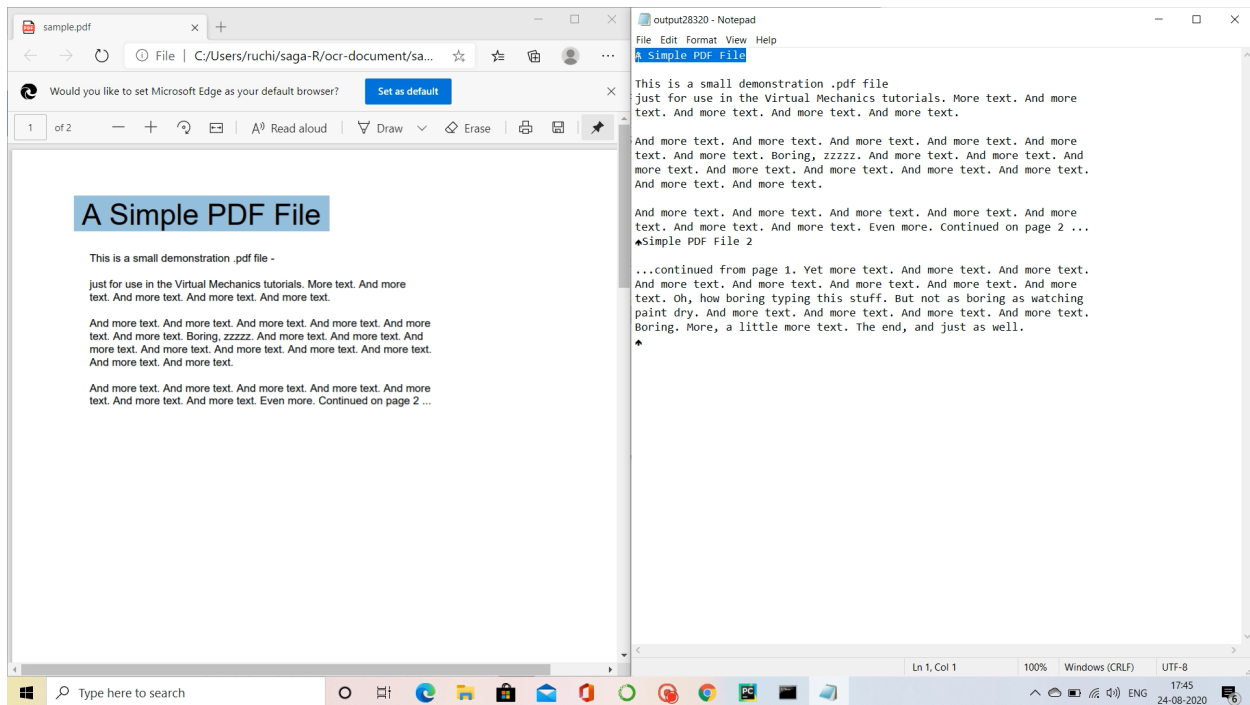
With the advent of OCR techniques, much time has been saved by automatically extracting the text out of a digital image of any invoice or a document.

Result:

Screenshots:







Applications:

- 1.Upload a pdf document, the document is analyzed by Optical character recognition package to extract text from it. The extracted text is again saved in a text document in the local drive.
- 2.Can be deployed in self driving cars,to read what is mentioned on road side safety board to get precaution while driving .

Conclusion:

Run python file on web,since it is deployed in flask,in the address you shown in command prompt.Upload a document file,ou will get a text file as output after processing through tesseract.

