

P.2. Text Extraction and Processing for Images, PDFs, and Speech Data

In text extraction applications, the core is the **Optical Character Recognition (OCR)** technology.

- Its primary function is to extract texts from images.
- Using advanced AI algorithms and **machine learning**, the OCR can identify and convert **image texts** into **audio files**, for easy listening.

There are some powerful text extraction software (having accuracy of 98+%), which are not freely/conveniently available. Thus we will implement two text extraction algorithms, one for image data and the other for speech data.

Project Objectives: The project will develop two separate Python programs: pdfim2text and speech2text.

1. **PDF-Image to Text** (pdfim2text)

- **Input:** an image or a pdf file
 - A PDF may include images.
 - A PDF may be generated by scanning, in which each page is an image.
- **Core Task:** Extract all texts; play texts from an image

2. **Speech to Text+Speech** (speech2text)

- **Input:** speech from **microphone** or a wave file
- **Core Task:** Extract texts; play the extracted texts.

An Example

```

pdfim2text
1  #!/usr/bin/python
2
3  import pytesseract
4  from pdf2image import convert_from_path
5  from PIL import Image
6  from gtts import gTTS
7  from playsound import playsound
8  import os, pathlib, glob
9
10 def takeInput():
11     pmode = 0;
12     IN = input("Enter a pdf or an image: ")
13     if os.path.isfile(IN):
14         path_stem = pathlib.Path(IN).stem
15         path_ext = pathlib.Path(IN).suffix
16         if path_ext.lower() == '.pdf': pmode=1
17     else:
18         exit()
19     return IN, path_stem, pmode
20
21 def pdf2txt(IN):
22     # you have to complete the function appropriately
23     return 'For pdf2txt, you may save the text here without return.'
24
25 def im2txt(IN):
26     # you have to complete the function appropriately
27     return 'For im2txt, try to return the text to play'
28
29 if __name__ == '__main__':
30     IN, path_stem, pmode = takeInput()    #pmode=0:image; pmode=1:pdf
31     if pmode:
32         txt = pdf2txt(IN)
33     else:
34         txt = im2txt(IN)
35
36     audio = gTTS(text=txt, lang="en", slow=False);
37     WAV = '0000_tmp.wav'; audio.save(WAV);
38     playsound(WAV); os.remove(WAV)

```

What to Do

First download **Image-Speech-Text-Processing.PY.tar**.

Untar it to see the file `pdfim2text` and some example codes in the directory `example-code`.

1. Complete `pdfim2text` appropriately.

- You may find clues from `example-code/pdf2txt.py`

2. Implement `speech2text` from the scratch.

- You may get hints from `speech_mic2wave.py` and `image2text.py` in the directory `example-code`.

Try to put all functions into a single file for each command, which enhances portability of the commands.

Report

- Work in a directory, of which the name begins with your last name.
- Use the three-page project document as an example for `pdfim2text`.
- zip or tar your work directory and submit via email.
- Write a report to explain what you have done, including images and wave files; upload it to Canvas.

PDF-Image to Texts

As a part of the project,
you will develop a Python program
that can extract texts from
PDF files and images:
and generally,
from PDF files including images.

An example PDF is
the one you are reading now.

(This portion is an image, by `text2image.py`.)

