

Continue



The script includes several classes and functions to accomplish this task. The `PDFExtractor` class is used to download and process PDF files. * The `download_file` function downloads a file from a specified URL. * The `break_pdf` method extracts individual pages from a PDF file and saves them as separate files. * There are two text extraction methods (`extract_text_algo 1` and `extract_text_algo 2`) that use different approaches to extract text from the PDF. The first method uses Python's built-in functionality, while the second method uses the PyPDF2 library with a custom device for conversion. * The `extract_text` function compares the extracted text using both methods and returns the longer one. * The `extract_table` function reads the table content from the PDF file using the tabula library. Additionally, there is a method called `tiff_header_for_CITIT` that generates a TIFF header for a given image size. `import PyPDF2 with open("sample.pdf", "rb") as pdf_file: read_pdf = PyPDF2.PdfFileReader(pdf_file) number_of_pages = read_pdf.getNumPages() page = read_pdf.pages[0] page_content = page.extractText() print(page_content)` Working with `LTTextBox`s from y-axis bottom up is feasible, allowing subtraction from the page's mediabox: `x0, y0, orig_x1, y1, orig_y0 = some_obj.bbox; y0 = page.mediabox[3] - y1, orig_y1 = page.mediabox[3] - y0, orig_y0`. Moreover, `LTTextBox`s possess a `.get_text()` method, which retrieves their textual content as a string. This is in addition to the `bbox` property. Each `LTTextBox` comprises `LTChars` (characters explicitly drawn by the PDF) and `LTAnnos` (spaces added by PDFMiner based on character distance). By combining these properties, you can display the coordinates of each text block. It's worth noting that `LTFigures`, unlike other Stack Overflow answers, do not require recursion since they directly contain `LTChar` objects rather than grouping them into `LTTextBox`s.

[Pdf to text remover](#). [Pdf to text to speech](#). [Pdf convert to text](#). [Pdf to text ai](#). [Pdf to text document](#). [Pdf to text ocr](#). [Pdf to text word](#). [Pdf to text editor free](#). [Pdf to text editor](#). [Pdf to text file](#). [Pdf image to text](#). [Pdf highlight to text](#). [Pdf convert to text free](#). [Pdf add to text](#). [Pdf to text copy](#).