# Final Report for EE447

Xu Jiayu 517030910367 and Hua Zeyu 517030910277

*Abstract*— Nowadays, we university students have a large demand for reading academic papers and 95 percent of these are stored in PDF files. Our group want to simplify the procedure of reading papers, so we design a tool that helps parsing the PDF file and extracting the information. In detail, we achieve such functions: We first parse the PDF file and get the page contents of it. Next we can extract the abstract, reference list from a paper, and then automatically download these references. Corely, we extract corresponding keywords cited by the original paper so that we can locate them quickly in the referenceto broaden the wanting knowledge instead of reading the whole reference blindly. Finally, we integrate all these functions and visualize them on a web page. All the codes have been uploaded to: **https://github.com/bilkerrr/EE447_FinalProject**.
Attention: in this report, the writer of each section also implements its function.

## I. INTRODUCTION (WRITTEN BY 517030910367 XU JIAYU)

We choose **PDF Parsing and Information Extraction** topic as our project. The motivation of this selection is that we are looking forward to improving the experience in reading academic papers stored in PDF file. We are always reading academic papers to broaden our minds and to keep pace with the technique development. However, tens-of-pages long and filled with quantities of reference, paper contents may make our lost in the paper. Therefore, we want to take avantage of PDF parsing tool to extract the key information of a paper so as to simplify the reading process and save time. In this project, we accomplish the following functions: (1) parse PDF file, (2) extract abstract of a paper, (3) extract reference list from a paper, (4) automatically downaload these references, (5) extract the key concepts of a reference cited by the original paper, (6) highlight these keywords in the reference PDF, (7) integrate all above functions and visualize them on a web page.

## II. PRELIMINARIES (WRITTEN BY 517030910367 XU JIAYU)

### A. PDF Background

Portable Document Format (PDF) is a file format developed by Adobe in the 1990s to present documents including texts and images. Based on the `PostScript` language, each PDF file encapsulates a complete description of a fixed-layout flat document whose structure is shown in 1, including the texts, fonts, vector graphics and other information. Howver PDF has no logical structures like sentences or paragraphs. Therefore, in order to get the real contents of a PDF we need to do further processing, called PDF parsing.



Fig. 1. Tree Structure of a PDF File

### B. PDFMiner Tool

PDFMiner[1] is a text extraction tool for PDF documents. It takes a strategy of lazy parsing, only to parse the stuff when needed, and hence it can save time and memory. To implement the parsing, it takes advantage of several classes shown in 2.

Fig. 2. Relationships between PDFMiner Classes

`PDFParse` fetches data from a file and `PDFDocument` stores it. Then `PDFPageInterpreter` to process the page cotents and `PDFDevice` translate it into the layout structure shown in 1. In our project, we can call the `LTTextBox` nodes to get the text information.

### C. RAKE Algorithm

Rapid Automatic Keyword Extraction (RAKE)[2] is proposed by Alyona Medelyan. The basic idea of the algorithm is that first to make phrase segmentation, score each word in the phrase, and rank the phrases with the sum score. To score each word, it first establishes co-present matrix for each word and calculates the degree in the matrix divided by the word frequency as its score. It works well in extracting jargons and runs fast.

## III. PDF PARSING (WRITTEN BY 517030910367 XU JIAYU)

In this part we use PDFMiner to get the text content of a PDF file. But though PDFMiner provides us with a convenient API, we still need further check and revise. For example, there might exist a link break in a word and a '-' will be attached to the end of a line such that, to restore the word, we need to remove the string `'-\n'`. What's more, some encoding problems need to be paid attention to like 'fi'.
The result is shown in 3.



Fig. 3. The Result of Abstract Extraction

## IV. INFORMATION EXTRACTION (WRITTEN BY 517030910367 XU JIAYU)

### A. Abstract Extraction

It is quite easy to implement. We only need to search for the independent paragraph `Abstract` and the next paragraph is what we are looking for.

### B. Reference List Extraction

Since reference are ordered in a special format, e.g. "[number] author name. paper title ...", we can take advantage of regular expression to fetch them. The result is shown in 4.



Fig. 4. The Result of Reference List Extraction

### C. Reference Downloading

Considering the patent problem, we only provide with downloading in ArXiv website. The implementation also makes use of regular expression to fetch the ArXiv ID of paper that recorded in the reference information. Then we can automatically generate the PDF URL according to the ArXiv ID for downloading.

## D. Keyword Extraction

We first need to do sentence tokenize and only preserve the sentences with citations. Sentence tokenize can be implemented according to the punctuations and Captial letters in the beginning of the sentence. Then we use RAKE algorithm to extract keywords of each sentence. After that, we match the keywords with corresponding reference. We split sentence into short clauses by puncuation marks, and then select the **nearest neighbour** citation number as the result of keyword matching. The intermediate keyword extraction result is shown in 5.
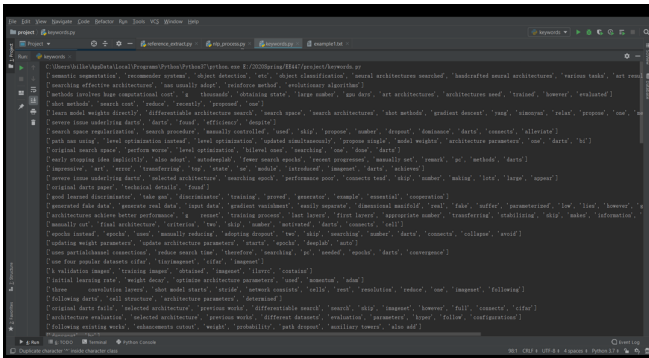


Fig. 5.   The Result of RAKE Algorithm

## E. Reference Highlighting

To search the keyword in the reference paper, we first need to process the keywords, because the phrase expression in each paper may vary. For example, the original paper says "gradient vanishment" but the reference uses the expression of "vanishing gradient". We need to extract stem of each keyword. Three common stemmers are widely used: PorterStemmer, LancasterStemmer and SnowballStemmer. After test, we choose SnowballStemmer due to its outperformance. Now, we want to highlight the keywords in the PDF. We use another PDF parsing tool: MuPDF[3]. It stands out among all similar products for its top rendering, and it can do more than PDFMiner. Here we use two functions: `searchFor` to locate the keywords and `addHighlightAnnot` to highlight them. A sample result is shown in 6. The original keywords are: **generator**, **real data** and **gradient vanishment**.



Fig. 6.   Highlighter Result

## V. VISUALIZATION (WRITTEN BY 517030910277 HUA ZEYU)

To visualize our work, we made a web application that is capable of extracting abstract and references as well as downloading highlighted referenced papers. User guide and how the application is implemented will be introduced in the following parts.

### A. User Guide

The homepage looks like a search engine, where you can upload a paper and submit to the server.



After the paper is uploaded, the extracted information is showed on the page. There is a 'BACK' button on the top of the page, redirecting to the home page for the next submission. The abstract is ranged under the button.

Scroll down the page and you will see the references. If the reference is from arXiv, there will be a 'FETCH PAPER' button under it. If you click the button, the corresponding referenced paper will be downloaded. Limited by copyright, however, we do not provide other access to papers, so non-arXiv papers cannot be downloaded from our application.



The last part is the highlighter. You can upload the referenced paper (you must renamed as its index first) and get a highlighted version. The highlighted keywords are extracted from the original paper you submit on the home page and that cites the current paper.



## B. Implementation

We deploy our project on apache server which is installed by xampp. The process of PDF analyzing and information extraction is written in Python and the web application is written in PHP, CSS, HTML and JavaScript.

*1) Call Python Scripts in PHP:* As is stated before, PDF analyzing and visualization parts are written in different languages, so we must find a way to connect Python to PHP.

The 'exec()' function can call Python scripts in PHP. It accepts three parameter, i.e. a string $command which is the command to be executed, an array $output that saves the result and an integer $return_var indicating the return value.

```php
$path="C:\‍\Python\python.exe ./pyfile/information_extract.py ";
$params = "C:/xampp/htdocs/source/".$fileInfoName;

exec($path.$params, $array, $res);
```
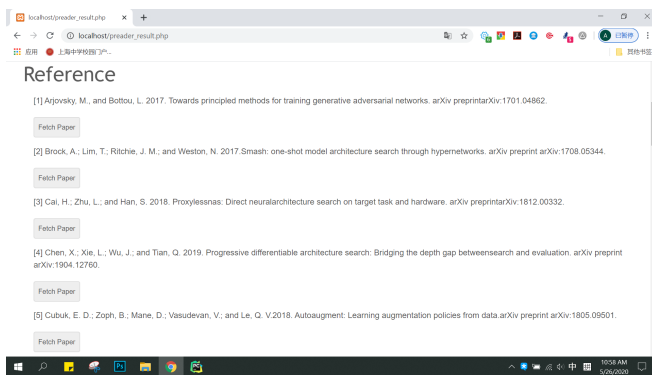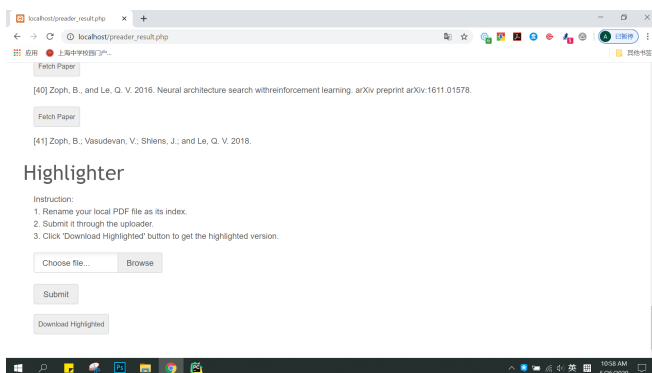
*2) Upload Files:* Uploading files is an important function in our application and we met most of the problems in this part.

**Size constraints** Since papers often contain pictures and tables, the size of PDF files is seldom less than 2MB. When we tried to upload files at first, we could successfully upload small files of hundreds of KBs, but we failed uploading files larger than 2MB. After searching some materials, we found out that this was because of the limits set in the configuration file of PHP. Here are some parmeters[5] to be revised:

- max_execution_time
- max_input_time
- post_max_size
- upload_max_filesize
- memory_limit

**Upload without refreshing** After we upload a paper on the home page, we are redirected to the second page. When we upload a referenced paper in the highlighter, however, we do not want to be redirected to another page, nor do we want to refresh the current page. However, since the file is uploaded through a form, there has to be a field of target.

The solution is to use iframe tag[4] in the page. iframe can embed a sub-html inside the large html. If we set target of the form as the iframe and when

we try to upload a paper in highlighter, only iframe will be refreshed, but others will not. Hence, it looks like the page is not refreshed at all.

```
echo '<div style="margin-left:30px">
    <div id="processing"></div>
    <form id="upload-form" action="upload.php" method="post" enctype="multipart/form-data"
        target="form-target" onsubmit="startUpload();">
        <input type="hidden" name="MAX_FILE_SIZE" value="5000000" />
        <label class="file">
        <input type="file" name="myfile" />
        <span class="file-custom"></span>
        </label>
        <br>
        <label class="submit">
        <input type="submit" name="submit" value="Upload" />
        <span class="submit-custom"></span>
        </label>
    </form>
    <br>
    <iframe style="width:0; height:0; border:0;" name="form-target"></iframe>
</div>';
```

## VI. FUTURE WORK (WRITTEN BY 517030910277 HUA ZEYU)

Although we have accomplished our goals in mid-term report, there is still a long way to go.

Firstly, we are currently using existed NLP model for general use, but we can train our own model to better fit keyword extraction in academic papers.

Secondly, we use Python for PDF extraction but PHP for visualization. As suggested by teaching assistant, we can use Django to constuct our website, which is written in Python and can better cope with PDF analyzing part.

Finally, we only extract keywords in referenced papers so far, but we can try to extract key sentences later.

## REFERENCES

[1] https://pypi.org/project/pdfminer/
[2] Rose, Stuart & Engel, Dave & Cramer, Nick & Cowley, Wendy. (2010). Automatic Keyword Extraction from Individual Documents. 10.1002/9780470689646.ch1.
[3] https://mupdf.com/
[4] https://www.cnblogs.com/wangmeijian/p/3978407.html?tdsourcetag=s_pctim
[5] https://blog.csdn.net/sinat_34328764/article/details/80053322?tdsourcetag=s_pctim_aiomsg