# *PDF Extraction and Clean-Up*

Sarah Bondurant, Nathan Davis, Richard Mays, Keegan Riley, Hayden Willeford
University of Arkansas, College of Engineering: Computer Science, Computer Engineering

## Introduction

This was a project given to us by Sorcero, an enterprise NLP suite created to support the Life Sciences & Insurance industries.

The arbitrary structure of a PDF is not useful for data extraction, and they would like to use Natural Language Processing, Computer Vision, and human intervention to train a model that would convert PDFs to document files.

## Purpose

The purpose of our project was to gain technical knowledge of PDFs, research possible conversion methods, and inspect existing implementations in order to find the best solution for Sorcero to pursue.
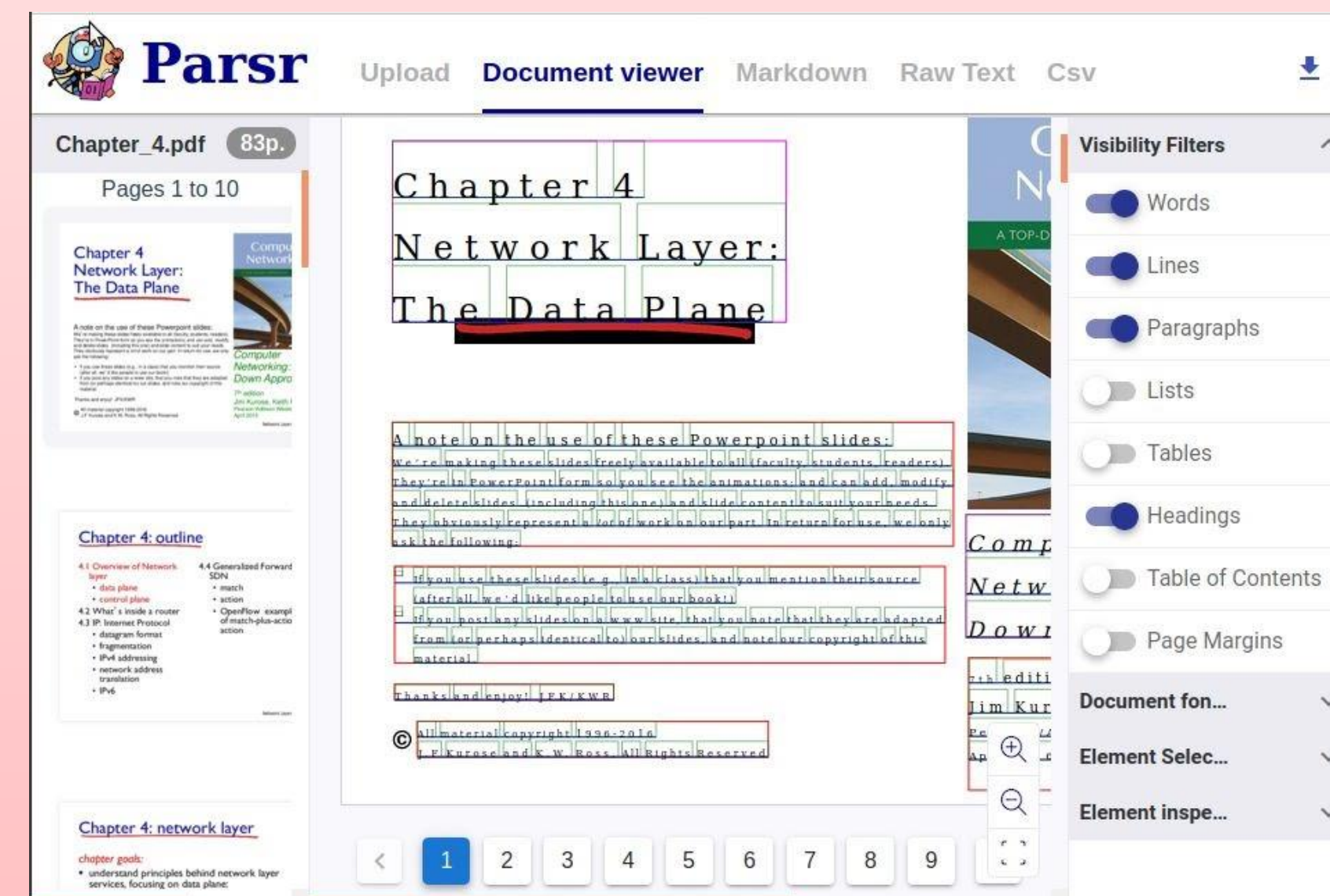
## Figure 2 (first image left)

Markdown output

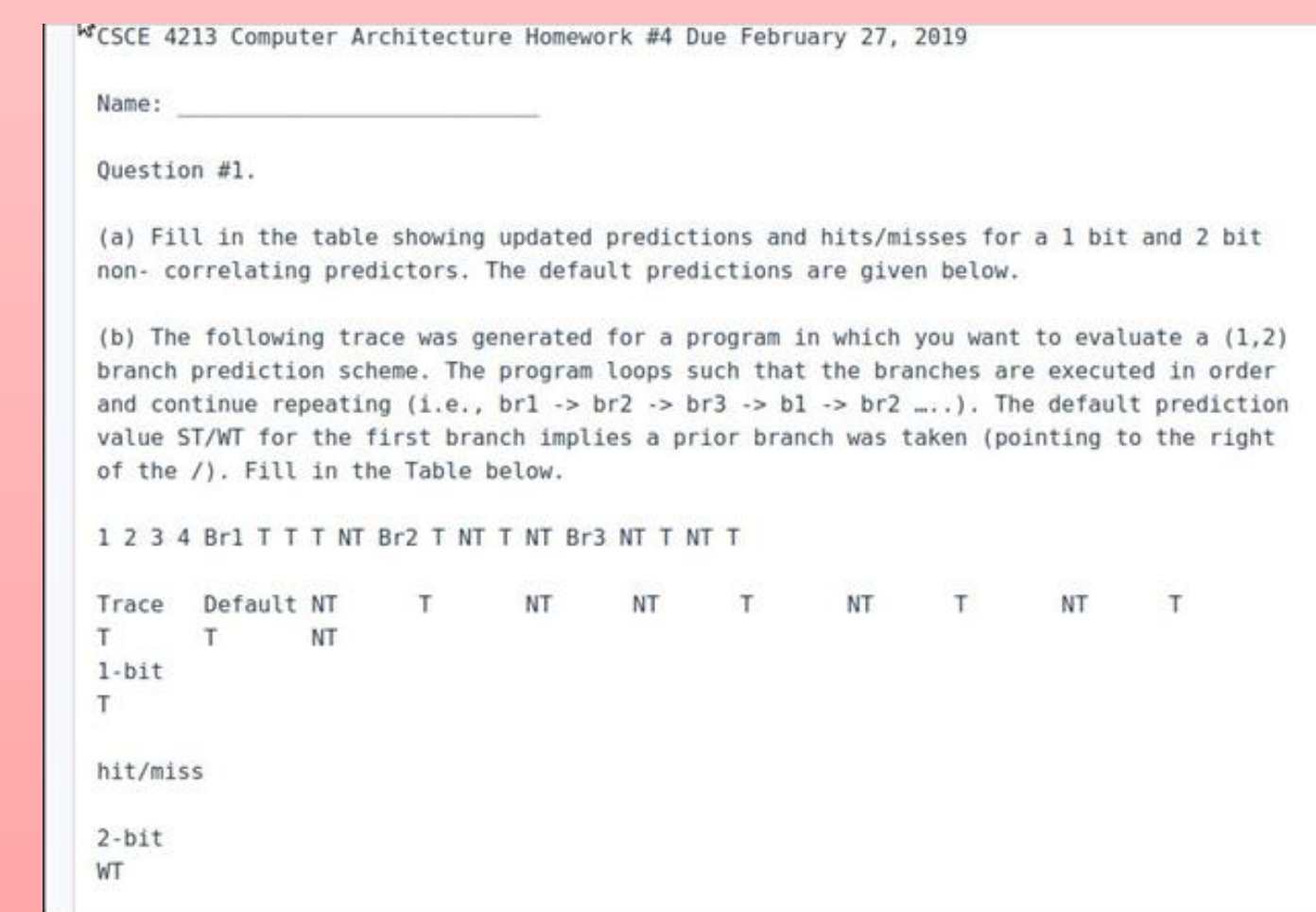## Figure 3 (second image left)

Raw text output

## Best Solution: Parsr

The best solution we found was Parsr, an open source document cleaning, parsing and extraction toolchain. Its strongest suit is the fact that it takes many already existing PDF parsers and combines them. However, Parsr is only available to run in a Linux Ubuntu environment with Docker installed. The program launches a GUI that can import PDF files for conversion.



## Figure 1 (image right)

The main output of Parsr includes the original document to the left, Parsr's output in the center, and the filters and element information to the right. Words are outlined in green, each line in blue, paragraphs in orange, and headers in pink.





## Other Solutions

XpdfReader is another open source PDF reader.

PyPDF4 is a python based PDF library.

Textricator is a tool that extracts text from documents.

Apache PDFBox is an open source Java tool for PDF extraction.

## Conclusion

When it comes to PDF translation, there are many methods that can be used in conjunction to get the best output. More research is needed to develop a solution for Sorcero's specific needs. However, Parsr fills very many of those needs. It combines the parsing capabilities of multiple different parsing programs to create a more powerful application.

The programs and methods we looked at in our research have their own pros and cons. Overall, the area of easy pdf translation that can accurately represent the original data is one that needs more research and development.