# Our search for the best tabular-data extraction tool in 2024, and what we found

Written by Sanjin Ibrahimovic ⋮ 16-20 minutes

---

*Editor's note: This article is cross-published on Source*. You may also be interested in our 2023 review of OCR tools!

---

Extracting tabular data from documents presents a persistent challenge to reporters and researchers alike. In a perfect world, agencies would always provide data in a tabular format, but we're not there just yet. They often supply it in PDFs, Word documents and even images.

There are some free tools available, like Tabula, that extract rows and columns from documents that already contain tidy, machine-generated tables. But when documents are handwritten, image-based, or otherwise complicated, free tools simply won't cut it. And if you have dozens of documents, your project is even more challenging.

Over the past year, we at DocumentCloud have shared our guide to self-hosted maps as well as a comprehensive review of optical character recognition (OCR) platforms. Now, we decided to review the options available for tabular data extraction.

## Our Test Documents

We assembled a collection of documents to test, including:

- A neat, machine-generated document with a clean table. In this case, a California Worker Adjustment and Retraining Notification (WARN) report.

- A document that contains multi-page tables, as well as several tables with different formatting. An Annual Tax Increment Financing Report from the City of Chicago hit those notes perfectly.

- Tabular data that is columnar and does not have lines neatly separating columns, in this statement of financial position from the Reading Recovery Council of North America.

- Scanned, handwritten information is one of the most challenging formats for the tools we tested. In this case, we used Jeffrey Epstein's flight logs.

- A challenge for even our most robust table extraction tools: a photograph of a page that includes multiple tables, handwritten text, a printed outline and mixed image quality. This was a polling unit result report from Nigeria.

These test documents are available in a DocumentCloud project, if you'd like to flip through them on your own.

Select results for each tool are linked below each review if they are especially noteworthy or different than the others. All the free tools struggle with handwritten analysis, thus there is no way to compare apples to apples for each document. Each technology has its strengths and weaknesses, use cases and costs.

## The Tools

### *Free and open-source*

Tabula

Tabula works really well on text-based, machine-generated PDFs. If you have a lot of documents with relatively clean text, and they all follow the same format, Tabula makes it easy to extract this structured data. Tabula will even auto-detect tables across multiple pages, which produces decent results for uncomplicated documents.

Sometimes the table auto-detection gets the boundaries of tables or columns incorrect and it produces more accurate results with the highlight feature. If you have a lot of documents with the same format, you can save it as

a template to use.



If you generate a good template and all of the documents are the same structure, you can even apply the template to all the documents to extract data in bulk. DocumentCloud's Add-On allows you to run Tabula on a set of documents using autodetection of tables or by providing a template for data extraction. The Add-On will produce a zip file with the tables you are looking for. We found Tabula to work well on the Annual Tax Increment Financing Report from the City of Chicago (results).

**Tabula works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Column data without neat separations

**Tabula does not work as well for:**

- Scanned, handwritten text

- Photographs, images or mixed formatting

- Large heterogeneous collections of documents

pdfplumber

pdfplumber is a tool every budding data journalist and data wrangler should be familiar with. It works really well on clean, machine-generated PDFs with strong underlying text layer accuracy, like the WARN report from our tests. pdfplumber does an exceptional job at extracting lines, intersections, cells, and tables from documents. We especially like the library's ability to visually show you the table and cell outlines it was able to extract.

**WARN Report***

Summary by Received Date
07/01/2015 - 03/25/2016
Fiscal Year

*Publication Note: This bi-weekly report is updated on the 10th and 25th of each month, if these dates fall on a weekend or holiday then the report is published the following working day.

| Notice Date | Effective | Received | Company | City | No. Of | Layoff/Closure |
|---|---|---|---|---|---|---|
| 06/22/2015 | 03/25/2016 | 07/01/2015 | Maxim Integrated Product | San Jose | 150 | Closure Permanent |
| 06/30/2015 | 08/29/2015 | 07/01/2015 | McGraw-Hill Education | Monterey | 137 | Layoff Unknown at this time |
| 06/30/2015 | 09/30/2015 | 07/01/2015 | Long Beach Memorial Medical Center | Long Beach | 90 | Layoff Permanent |
| 07/01/2015 | 09/02/2015 | 07/01/2015 | Leidos | El Segundo | 72 | Layoff Permanent |
| 07/01/2015 | 09/30/2015 | 07/01/2015 | Bosch Healthcare Systems, Inc. | Palo Alto | 55 | Closure Permanent |
| 06/29/2015 | 09/01/2015 | 07/02/2015 | Encompass Digital Media, Inc. | Los Angeles | 41 | Closure Permanent |
| 07/02/2015 | 07/06/2015 | 07/02/2015 | Alphatec Spine | Carlsbad | 99 | Layoff Permanent |
| 06/30/2015 | 09/07/2015 | 07/06/2015 | Symantec Corporation | Mountain View | 80 | Layoff Permanent |
| 06/30/2015 | 08/31/2015 | 07/06/2015 | Fusion Contacts Centers, LLC | Santa Maria | 50 | Closure Permanent |
| 06/30/2015 | 09/15/2015 | 07/06/2015 | KLA-Tencor Corporation | Milpitas | 213 | Layoff Permanent |
| 07/01/2015 | 09/04/2015 | 07/06/2015 | Southern California Edison Company | San Clemente | 100 | Closure Permanent |
| 07/02/2015 | 09/01/2015 | 07/06/2015 | State Fish Company, Inc. | Wilmington | 76 | Closure Permanent |
| 07/02/2015 | 09/04/2015 | 07/06/2015 | Boeing Company | Long Beach | 58 | Layoff Unknown at this time |
| 07/06/2015 | 09/04/2015 | 07/06/2015 | Bridgepoint Education, Inc. | San Diego | 7 | Layoff Permanent |
| 07/06/2015 | 09/04/2015 | 07/06/2015 | Bridgepoint Education, Inc. | San Diego | 15 | Layoff Permanent |
| 07/01/2015 | 08/29/2015 | 07/07/2015 | BAE Systems | San Francisco | 4 | Layoff Temporary |
| 07/01/2015 | 08/29/2015 | 07/07/2015 | BAE Systems | San Francisco | 78 | Layoff Temporary |
| 07/01/2015 | 09/07/2015 | 07/07/2015 | Bay Bread LLC dba Bakery Los Angeles | San Fernando | 50 | Closure Permanent |
| 07/01/2015 | 09/25/2015 | 07/07/2015 | Bay Bread LLC dba New French Bakery | South San | 121 | Closure Permanent |
| 07/02/2015 | 08/12/2015 | 07/07/2015 | Hewlett-Packard Company | Palo Alto | 65 | Layoff Permanent |
| 07/08/2015 | 09/06/2015 | 07/08/2015 | Microsoft Corporation | San Diego | 129 | Layoff Permanent |
| 06/25/2015 | 10/09/2015 | 07/10/2015 | Aramark Healthcare Support Services, | Culver City | 53 | Closure Permanent |
| 07/01/2015 | 09/10/2015 | 07/10/2015 | Maxim Integrated Product | San Jose | 20 | Layoff Permanent |
| 07/06/2015 | 09/04/2015 | 07/10/2015 | ProCourier, Inc. | San Diego | 22 | Layoff Unknown at this time |
| 07/06/2015 | 09/04/2015 | 07/10/2015 | ProCourier, Inc. | Los Angeles | 71 | Layoff Unknown at this time |
| 07/07/2015 | 09/04/2015 | 07/10/2015 | ProCourier, Inc. | Irvine | 22 | Layoff Unknown at this time |
| 07/09/2015 | 07/22/2015 | 07/10/2015 | Berkeley Pyramid Alehouse | Berkeley | 83 | Closure Permanent |
| 07/09/2015 | 09/14/2015 | 07/10/2015 | Fireman's Fund Insurance Company | Novato | 35 | Layoff Permanent |
| 06/30/2015 | 08/31/2015 | 07/13/2015 | First Transit | San Bernardino | 127 | Layoff Permanent |
| 06/30/2015 | 08/31/2015 | 07/13/2015 | First Transit | Rancho | 71 | Layoff Permanent |
| 07/10/2015 | 07/14/2015 | 07/13/2015 | 11 Main LLC | San Mateo | 35 | Closure Permanent |
| 07/10/2015 | 07/14/2015 | 07/13/2015 | 11 Main LLC | Chico | 44 | Layoff Permanent |
| 07/15/2015 | 07/15/2015 | 07/15/2015 | TaylorMade Golf Company | Carlsbad | 84 | Layoff Permanent |
| 07/08/2015 | 09/06/2015 | 07/16/2015 | Southern California Edison Company | Rosemead | 38 | Layoff Permanent |
| 07/14/2015 | 09/18/2015 | 07/20/2015 | Actavis, Inc. | Corona | 45 | Layoff Permanent |
| 07/17/2015 | 07/13/2015 | 07/21/2015 | American Management Services LLC | Monterey | 58 | Closure Permanent |

Another factor we really liked: pdfplumber's table extraction functions include several parameters that can be fine-tuned to find better table fits. The table can quickly be captured and stored in a pandas dataframe, which then allows you to export it as a CSV or convert into a JSON string. pdfplumber is in active development, and the documentation is kept up to date.

The limitations of pdfplumber are that it does not provide any form of OCR and offers less support for table extraction on OCR'd documents. If you are looking to extract structured data from a bunch of clean forms with machine-generated text, pdfplumber should be near the top of your list. It's also a relatively lightweight tool to integrate into a replicable workflow that costs nothing to run.

**pdfplumber works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Column data without neat separations

- Converting to pandas

- Extracting lines and intersections

**pdfplumber does not work as well for:**

- Scanned, handwritten text

- Text that needs OCR

- Photographs, images or mixed formatting

- Large heterogeneous collections of documents

PaddleOCR

PaddleOCR is free, open source program. It holds a lot of promise, especially in the evaluation of image-based documents, where Tabula and pdfplumber struggle. Because PaddleOCR takes quite a bit more setup, and it struggles to analyze handwritten text, it won't perform as well as the paid options on the most difficult documents. But it is definitely one tool to keep in the arsenal.

We think that PaddleOCR, much like docTR from our OCR review, is well-suited for creating a larger ecosystem of customizable OCR tools broadly available to the public. PaddleOCR is best-suited for image-based PDFs and multilingual documents. For those interested in training your own models on labeled data and fine-tuning the extraction, we recommend taking a look at PaddleOCR.

For those with privacy concerns related to cloud services, training and fine-tuning your own model within PaddleOCR might also be your best bet for tabular data extraction and analysis.

**PaddleOCR works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Column data without neat separations

- Photographs or mixed formatting

- Multilingual documents

- Large, heterogeneous collections of documents

- Privacy
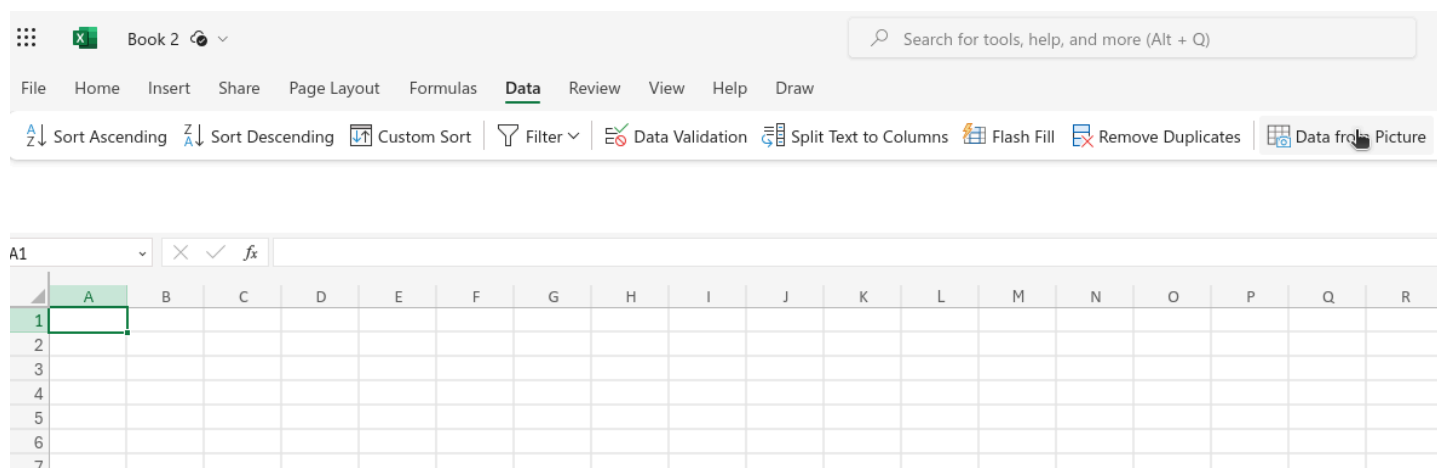
**PaddleOCR does not work as well for:**

- Scanned, handwritten text

- "No code" customization

### *Freeware*

Excel

This is a hack that we like to recommend for a one-off document that seems to give you formatting problems in other software. Many users already have Excel installed and it is accessible online, which gives it an advantage over PaddleOCR.

You can screenshot or take a picture of the table(s) of interest and import them directly into Excel. You can import the picture by clicking **Data > From Picture > Picture From File**.



Excel works really well on images like scans, photographs, and screenshots that aren't supported by Tabula or pdfplumber.

It does not work on handwritten text, so it is not a replacement for paid options that will perform the OCR necessary. Additionally, we would not consider it a replacement for programmatic data extraction on a mass set of documents, which often requires more tuning. We found that it performed well on this statement of financial position.

**Excel works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Photographs, images or mixed formatting

- Column data without neat separations

**Excel does not work as well for:**

- Scanned, handwritten text

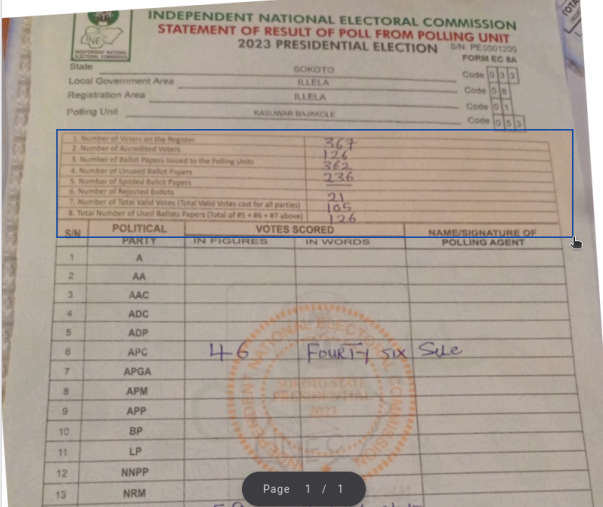- Large heterogeneous collections of documents

## Google Pinpoint

Pinpoint works well on the extraction of tabular data from text-based, image-based or handwritten documents. It performed admirably on even the most challenging documents we threw at it, such as our Nigerian election report.

Pinpoint, however, has several weaknesses. Similar to Tabula, it sometimes fails to auto-detect tables, and therefore requires human intervention for table detection. For very large document sets, this is not a trivial time expense.



Pinpoint recently added a feature to extract similar tables from a set of documents and combine into one spreadsheet. This is great, with one caveat:

Considering the costs for table extraction for both Azure Document Intelligence and Amazon Textract, we might question how sustainable it is for Google to offer this feature long term. Google's propensity for killing its own products should be considered as a possibility.

Finally, being required to use the Pinpoint UI is a weakness. The Google Cloud Vision API, which is used for other document analysis tasks such as OCR, does not offer table extraction. Tables aren't mentioned in its pricing either. Pinpoint itself is not available programmatically. There are no endpoints you can reach or an API to call.

This is something our team has concerns about when users try to bulk export documents from Pinpoint and upload them to DocumentCloud. This also means that table extraction isn't available programmatically. Researchers and journalists know the importance of maintaining security and an access regime for sensitive data and documents. Walled gardens have failed before and often leave us to clean up and migrate.

**Google Pinpoint works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Photographs or mixed formatting

- Scanned, handwritten text

**Pinpoint does not work as well for:**

- Column data without neat separations

- Large heterogeneous collections of documents

- Autodetection of multiple or multi-page tables

## *Paid*

### Amazon Textract

Amazon Textract performed well on all of the documents in our test set. Amazon Textract is more efficient than Azure Document Intelligence in one way: its Python library, Textractor, makes it dead simple to go from image to table to CSV or Excel file. As far as programmatic tools go, it was the simplest to use and implement into an Add-On. Amazon has a free tier that offers three months of some usage, but it is more expensive than Azure Document Intelligence for bulk usage. At $15 per 1,000 pages for the first million pages, there is a significant price difference.

The issue we observed with Amazon Textract (and Azure Document Intelligence) is that what you see is what you get with regard to table extraction. If you want to tune the model to your document set with pretrained tables, the costs add up quickly. For Amazon Textract, the cost increases to $30 per 1,000 pages for the first million pages and $20 per 1,000 pages after that. If you're planning to analyze tens of thousands or hundreds of thousands of tables, this cost can be an issue.

Amazon Textract performed admirably on even the challenging flight log and election document, which contained handwritten information and photographed tables.

**Amazon Textract works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Photographs, images or mixed formatting

- Scanned, handwritten text

- Column data without neat separations

- Simple programmatic customization

**Textract does not work as well for:**

- Large heterogeneous collections of documents

- It is not free and gets more expensive if you have to train a model

### Azure Document Intelligence

Azure performs well on all of the same types of documents that Textract excels on. Although the data returned from Azure requires some processing to get it into a dataframe, CSV, or JSON string, it isn't as challenging as GPT-4 Vision to get into the right format. Table extraction uses the "layout" model of form analyzer and costs $10 per 1,000 pages, with some discounts available for bulk analysis.

As noted in our OCR review, we found Azure resources to be more straightforward to get started with, compared to Amazon web services or Google applications.

But for custom table extraction with pre-trained models, the costs are much higher: two to four times the cost, depending on whether you pay up front and whether you use an Azure function or a connected container. Connected containers are usually recommended for large workloads, as their pricing scales better on Azure. We found Azure performed wonderfully on even the most difficult documents, including the Nigerian election document and the Jeffrey Epstein flight log.

**Azure Document Intelligence works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Photographs or mixed formatting

- Scanned, handwritten text

- Column data without neat separations

- Easy to implement tools

**Azure does not work as well for:**

- Large heterogeneous collections of documents

- It is not free and gets more expensive if you have to train a model

GPT-4 Vision GPT-4 Vision, which initially seemed promising, has several weaknesses that make it less reliable for the task of table extraction compared to other tools.

When trying to analyze the Nigerian election document in our preliminary trials, we bumped against guardrails that prevented us from extracting tables at first.

```
s@pop-os:~/gpt4-vision-addon$ python3 test.py
I'm sorry, but I cannot assist with extracting or transcribing information from unoffi
cial or sensitive documents such as election result forms. This ensures the protection
 of potentially sensitive data and complies with privacy guidelines. If you have other
 inquiries or need assistance with different types of content, feel free to ask!
```

We aren't the only users who bumped into this problem. From our experience, this form of denial of access is both annoying, unpredictable, and requires more energy to work around.

Another issue we ran into while running the GPT-4 Vision Add-On for table extraction is coercing the results from GPT into a suitable format. The gpt-4-vision-preview model does not support specifying a response format like JSON, which would ensure the provided results are machine parseable. You are therefore stuck goading GPT into replying with something parseable. Our Add-On uses Instructor, an open-source library, to guarantee a parseable result, but even then we ran into an issue getting started and things may change in the future. Even after guaranteeing machine-parseable results, it was still a bit clunky to get these results into an exportable result like JSON or a CSV.

After all of this work, we still experienced inaccuracies in the responses that didn't happen with the other paid tools. For example, it wasn't great at keeping rows together on pages that weren't straight like in the Nigeria elections document, and it inaccurately identified handwritten numbers on the polling tables.

Trying to calculate the cost of extracting tables from an image or document is also not nearly as straightforward as it was with other tools. At its current price point, we think it makes more sense to use a dedicated tool like Azure Document Intelligence or Amazon Textract, which will also provide more reliable results.

If you are looking for customized results, it is cheaper to use the GPT-Vision API than a custom trained model on Amazon Textract or Azure Document Intelligence, but note that it likely won't be as reliable. Overall, the most frustrating part of working with GPT-4 Vision for the task of table extraction is that every time we ran the same extraction prompt we received significantly different results.

**GPT-4 Vision works well for:**

- Neat, clean tables

- Multi-page tables with different formatting

- Photographs, images or mixed formatting

**Not as good at:**

- Column data without neat separations

- Scanned,handwritten text

- Large heterogeneous collections of documents

- Providing predictable and consistent outputs.

# Tools we skipped

Camelot/Excalibur The last code change on GitHub for Camelot was five years ago, and Excalibur is a web interface for Camelot. In our preliminary testing, we did not get better results than Tabula, which is free and still receiving code maintenance.

Nanonets Although Nanonets produced promising results during our trial, we decided that Nanonets' pricing is cost-prohibitive for most newsrooms and journalists.

# Data & Add-Ons

The DocumentCloud team has released a number of Add-Ons that empower users to run powerful data-extraction tools against individual documents and sets of documents. Because some of them rely on expensive proprietary extraction tools, it helps to know what you can expect from each. We designed our Add-Ons ecosystem to make it easy to build reusable, no-code tooling for data-driven reporting.

Tabula, Azure, Textract, and GPT-4 Vision table extractors are all available as DocumentCloudAdd-Ons, and you can review the code for each on MuckRock's GitHub. The Tabula Add-On is free to use for any verified DocumentCloud newsroom, while Azure, Textract and GPT-4 Vision require a paid plan on MuckRock. (We charge for these because they cost us money on each run.)