# Assignment 2

Drive Link - https://drive.google.com/drive/folders/1VozeLR1GTBusVJXih_fo1F7lMnY-L-HG?usp=drive_link

Codelabs Link - https://codelabs-preview.appspot.com/?file_id=1-5QP7m-QK3vR2Jtv8-lNfvdp06lYvoHcKmNlVWzhRqU#0

Deployment link - https://18.188.86.27:8080

# PDF Extraction API Evaluation Template

**Team: Team 5**

**Team members: Saurabh Vyawahare, Aniket Patole, Shreya Bage**

**Link to your analysis:**

**Summary:**

The evaluation compared two methods of PDF text extraction: **OpenAI's GPT-based API** and **PyPDF**. OpenAI's API offers advanced document understanding, is scalable for large volumes, and supports multi-language extraction. It excels at handling complex documents and providing context-aware responses but comes with higher costs due to usage-based pricing. PyPDF is more suited for small-scale, local extraction tasks, with limited AI features but lower costs. Both solutions have pros and cons, making PyPDF ideal for straightforward extraction needs, while OpenAI is better for large-scale, complex, or multi-language extractions.

## 1. General Information

| Attribute | Details |
|---|---|
| **API Name** | PyPDF and OpenAI Extractor |
| **Vendor** | PyPDF: Open-source library.<br>OpenAI: OpenAI as the vendor. |
| **Version/Release Date** | PyPDF: PyPDF 5.0.1<br><br>OpenAI: GPT 4-o |
| **Pricing Model** | PyPDF: Free.<br>OpenAI: Pay-as-you-go based on API usage (tokens). |

| Licensing and Compliance | PyPDF: Open-source licenses (e.g., MIT). OpenAI: GDPR, CCPA, |
|---|---|

## 2. Technical Capabilities

| Feature | Adobe PDF API | OpenAI API | PyPDF | Other Vendor |
|---|---|---|---|---|
| **File Format Support (PDF, DOCX, etc.)** | Supported formats for input/output. | Can handle various file formats indirectly via extraction methods. | PDF (input/output). | |
| **OCR (Optical Character Recognition)** | Accuracy of OCR for extracting text from images. | Can extract text, but not specifically designed for OCR. | Does not natively support OCR. | |
| **Table Extraction** | Ability to accurately detect and extract tables. | May struggle with complex table extraction but can understand structured data better with context. | Limited to parsing tables as plain text. | |
| **Form Extraction** | Handling of structured forms (checkboxes, radio buttons). | May be able to interpret structured text with proper prompting. | Cannot extract structured forms. | |
| **Complex Layout Support** | Ability to extract data from complex layouts (columns, images, embedded objects). | Can handle complex layouts better due to its language model capabilities. | Handles simple layouts. | |
| **Multi-language Support** | Support for extracting content in multiple languages. | Supports multiple languages natively. | Supports only the character sets in PDFs. | |

| | | | | |
|---|---|---|---|---|
| **Scalability and Performance** | Speed and ability to scale across large datasets. | Scales based on OpenAI infrastructure. | Limited by your processing power. | |
| **API Integration and Usability** | Ease of API integration, SDK availability, documentation quality. | Easy to integrate via API but requires token management and optimization. | Easy integration for basic extraction. | |
| **Customization Options** | Customizable extraction rules, fine-tuning. | Highly customizable via prompts. | Limited customization. | |
| **Accuracy and Error Handling** | Metrics on accuracy, handling of errors and ambiguous data. | More accurate for understanding, but results can vary depending on prompt design. | Simple extraction with frequent errors in complex documents. | |

## 3. Business and Strategic Considerations

| Evaluation Metric | Adobe PDF API | OpenAI API | PyPDF | Other Vendor |
|---|---|---|---|---|
| **Cost Efficiency (Pricing vs. Features)** | Balance between cost and the number of features offered. | Costs increase with usage. | Free | |
| **Vendor Reputation and Stability** | Market position, experience, and reliability of vendor. | High reputation in AI and language processing. | Community-driven. | |
| **Customer Support and SLA** | Quality and response times of customer service, availability of SLA. | Has dedicated customer support and documentation. | Community-based support. | |
| **Security and Privacy** | How the API handles data encryption, anonymization, and compliance with regulations (GDPR, CCPA). | Follows security standards but involves sending data to external servers. | Depends on your security setup. | |
| **Documentation and Training Resources** | Availability of user guides, developer resources, and training materials. | Extensive documentation with many use-case examples. | Well-documented but limited examples. | |
| **Community and Ecosystem** | Size and engagement of user community, availability of third-party integrations, plugins. | Large, highly engaged community, extensive third-party integrations. | Medium-sized open-source community, limited third-party integrations. | |
| **Roadmap and Innovation** | The vendor's future plans for the API, commitment to innovation. | Fast-paced innovation, regular updates, advanced AI research. | Community-driven, incremental updates, limited innovation. | |

| Vendor Lock-in Risk | Risk of dependency on the API, ease of migration to alternative solutions. | Moderate to high dependency, more complex migration to other AI platforms. | Low dependency, easy migration to other open-source tools. | |
|---|---|---|---|---|

## 4. Performance Metrics

| Metric | Adobe PDF API | OpenAI API | PyPDF | Other Vendor |
|---|---|---|---|---|
| Latency | Average response time for extraction requests. | Slower due to external API calls. | Fast as it's processed locally. | |
| Throughput | Number of pages processed per minute/hour. | Can process larger datasets via scalable API. | Dependent on local resources. | |
| Error Rate | Frequency and types of errors during processing. | Fewer errors with well-constructed prompts but can misinterpret complex structures. | Higher error rates in structured data. | |
| Data Loss/Integrity | Percentage of lost/misinterpreted data during extraction. | Low data loss, high accuracy. | Moderate data loss with complex layouts. | |

## 5. Value-Add Features

| Feature | Adobe PDF API | OpenAI API | PyPDF | Other Vendor |
|---|---|---|---|---|
| **Advanced AI/ML Capabilities** | AI/ML support for better contextual extraction, document understanding. | Strong ML/AI capabilities for understanding text context. | None. | |
| **Pre-built Templates for Specific Use Cases** | Availability of industry-specific extraction templates (e.g., legal, financial). | No direct templates, but you can create custom prompts. | None. | |
| **Document Classification/Tagging** | Auto-classification of documents based on content. | Can auto-classify documents based on context using AI. | No built-in auto-classification capabilities. | |
| **Metadata Extraction** | Ability to extract embedded metadata (author, timestamp, etc.). | Can extract and analyze embedded metadata in documents. | Basic metadata extraction (author, timestamps, etc.). | |

## 6. Overall Evaluation

| Attribute | OpenAI Rating | Comments |
|---|---|---|
| Technical Fit | 8/10 | best for complex document understanding |
| Business Fit | 8/10 | scalable, but costs can rise |
| Total Cost of Ownership | 7/10 | usage-based costs |
| Ease of Implementation and Use | 8/10 | requires API management |
| Vendor Reliability and Support | 9/10 | Reliable with frequent updates and good community support. |

## 7. Recommendations

| Recommendation | Details |
|---|---|
| Best Fit for the Use Case | PyPDF: Best for small-scale, local PDF processing.<br>OpenAI: Best for advanced document understanding, large-scale processing, and multi-language support. |
| Further Considerations | Consider potential costs for scaling with OpenAI and data privacy concerns for cloud-based solutions. |