

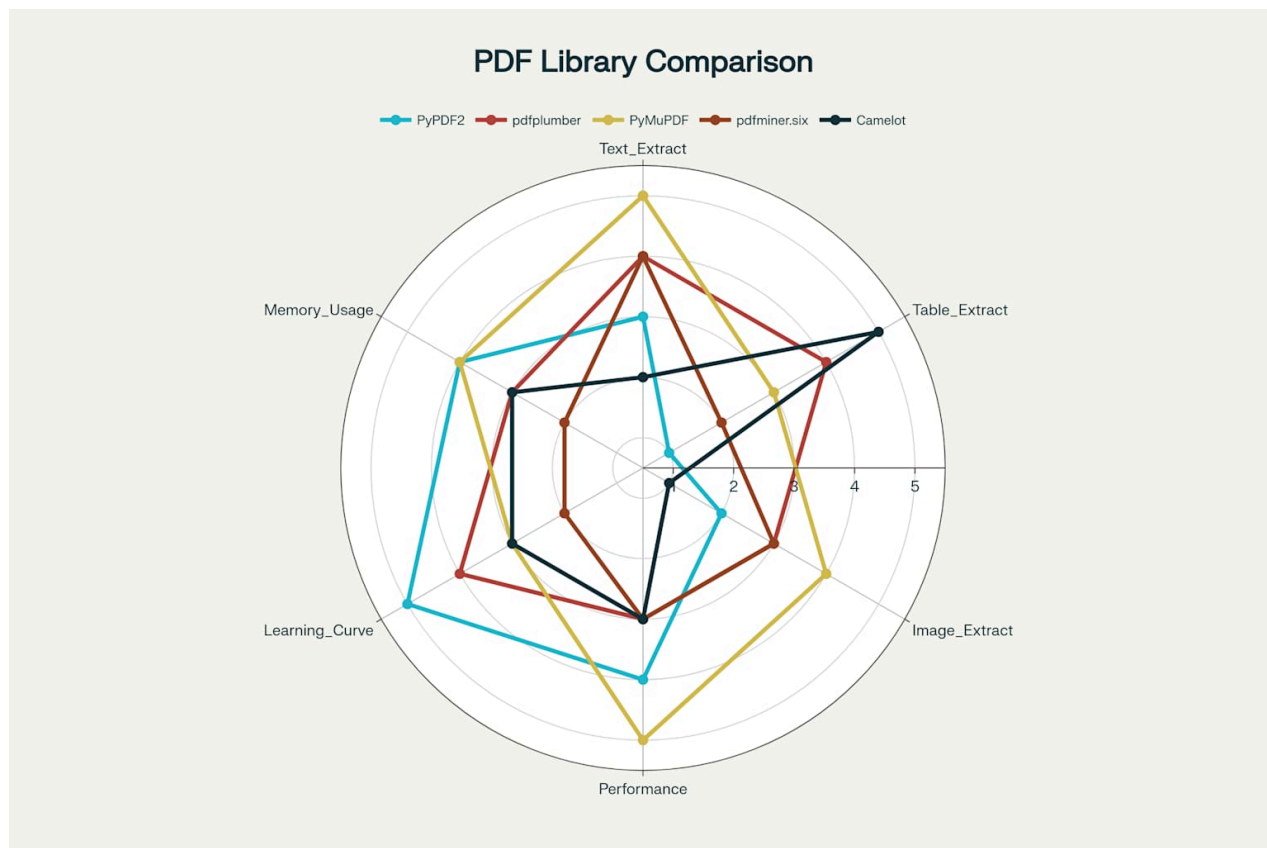
# PDF 資料提取 Web API 開發完整指南

隨著數位化轉型的深入發展，PDF 文件作為企業和學術機構的主要資料載體，其自動化資料提取需求日益增長。本報告基於最新研究成果和實務經驗，為開發高效能、安全可靠的 PDF 資料提取 Web API 提供全面的技術指導，涵蓋從底層解析技術選擇到高層架構設計的各個關鍵面向。

## PDF 解析技術基礎與工具選擇

### 主流 PDF 解析庫比較分析

現代 PDF 解析生態系統提供了多種技術方案，每種都有其特定的優勢和應用場景。根據最新的比較研究，不同解析工具在文字提取、表格識別、圖像處理等方面表現差異顯著<sup>[1]</sup>。PyMuPDF 和 pypdfium 在文字提取方面通常表現最佳，但所有解析器在處理科學文獻和專利文件時都面臨挑戰<sup>[1]</sup>。對於這些困難類別，基於學習的工具如 Nougat 展現出卓越的性能<sup>[1]</sup>。



### Comprehensive comparison of PDF parsing libraries across key capabilities

在表格檢測方面，Table Transformer(TATR) 在金融、專利、法律法規和科學類別中表現優異<sup>[1]</sup>。針對招標文件，Camelot 工具表現最佳，而 PyMuPDF 在手冊類別中表現卓越<sup>[1]</sup>。這些發現強調了根據文件類型和特定任務選擇適當解析工具的重要性<sup>[1]</sup>。

對於複雜的 PDF 文件處理，研究顯示結合多個庫可以有效應對不同挑戰<sup>[2]</sup>。例如，可以使用 pdfplumber 處理表格，PyPDF2 處理文字提取，當文件同時包含結構化和非結構化內容時<sup>[2]</sup>。對於掃描的 PDF，建議結合 pdf2image 和 Tesseract OCR 來處理基於圖像的文字提取<sup>[2]</sup>。

### 基於人工智慧的提取技術

大型語言模型（LLM）在證據合成的資料提取方面展現出巨大潛力。比較研究顯示，Claude 2 的準確率高達 96.3%，而使用第三方插件的 GPT-4 準確率為 68.8%，但大部分錯誤來自插件本身<sup>[3]</sup>。當提供 PDF 的選定文字時，Claude 2 和 GPT-4 分別準確提取了 98.7% 和 100% 的資料元素<sup>[3]</sup>。

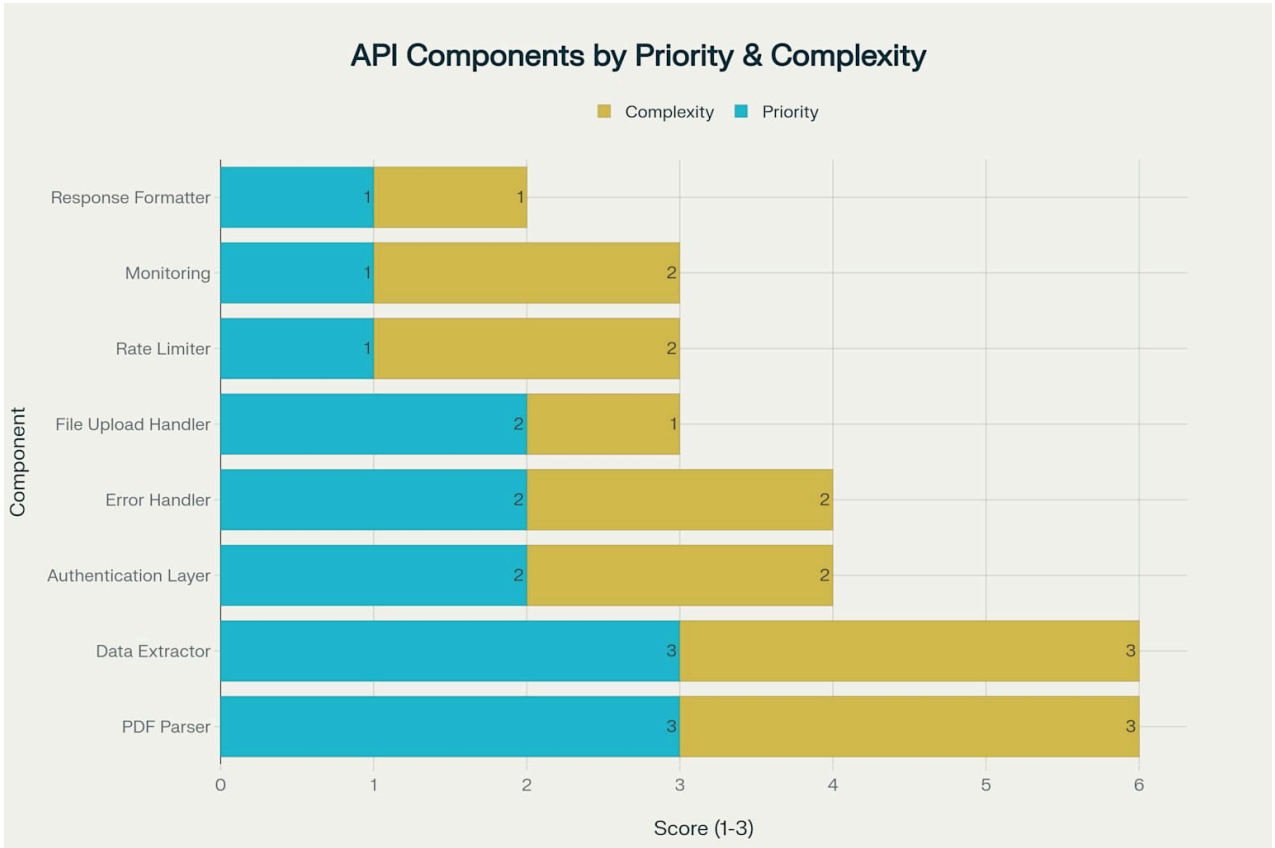
Adobe PDF Extract API 利用 Adobe Sensei AI 技術提供高度準確的資料提取，適用於原生和掃描 PDF，無需自訂機器學習模板或模型訓練<sup>[4]</sup><sup>[5]</sup>。該服務能夠提取文字、複雜表格和圖形，並以結構化 JSON 格式輸出內容<sup>[4]</sup>。

### API 架構設計與實作策略

#### 核心架構組件

開發 PDF 資料提取 Web API 需要考慮多個關鍵組件的協調配合。從優先級和複雜度角度分析，核心組件包括認證層、檔案上傳處理器、PDF 解析器、資料提取器等

。



API Architecture Components ranked by Priority and Implementation Complexity

RESTful API 設計應遵循 2024 年的最佳實踐，包括基於資源的架構、無狀態通訊、客戶端-伺服器分離、統一介面等<sup>[6]</sup>。關鍵建議包括使用清晰的名詞基礎資源名稱、實施 OAuth 2.0 和速率限制、快取

資料和壓縮回應、提供 OpenAPI 互動式文件、使用語義版本控制等 [6]。

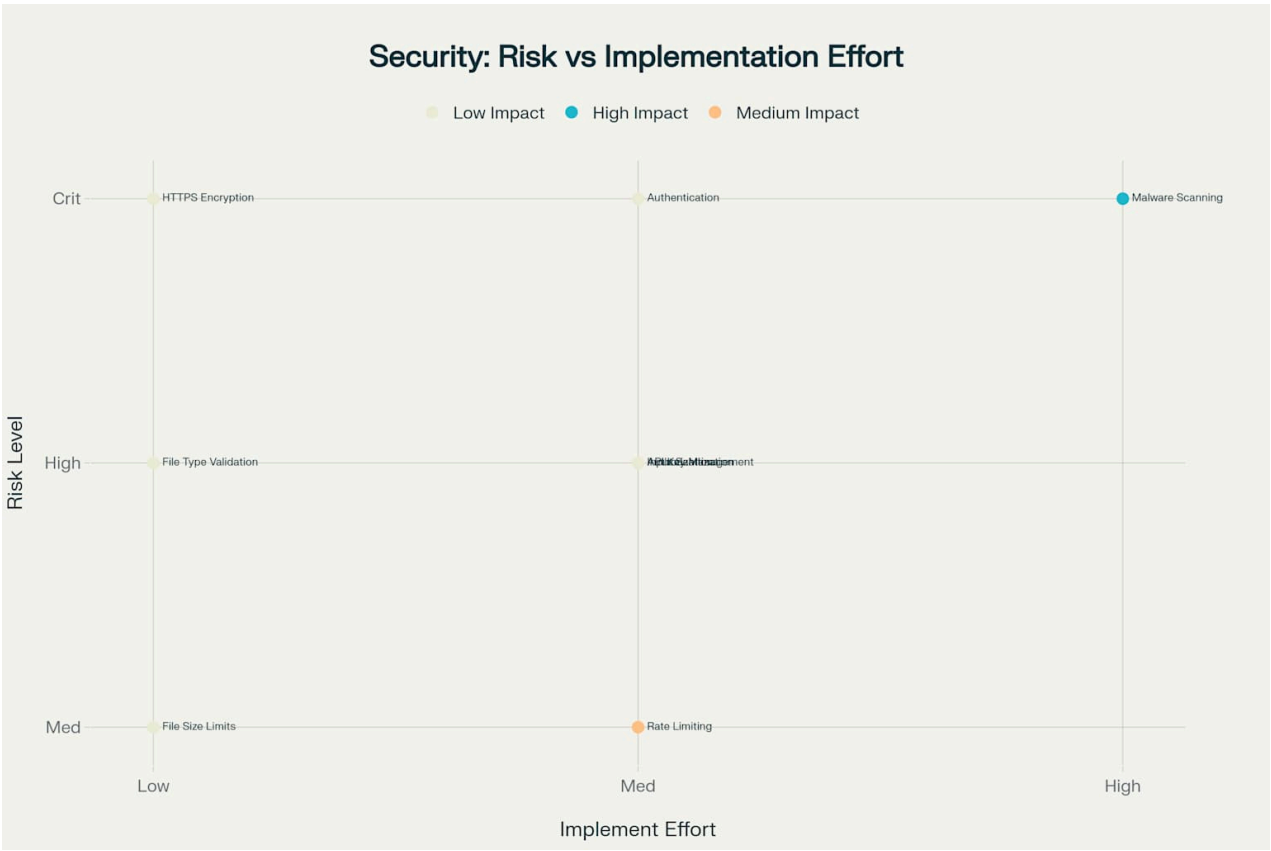
### 檔案上傳處理機制

處理檔案上傳時，應始終使用 multipart/form-data 內容類型和 POST 請求方法 [7]。這是因為檔案通常以二進位資料發送，需要特定的編碼方案確保正確傳輸 [7]。對於大型檔案，建議實施分塊上傳，將大檔案分解為較小的塊並非同步處理以提高性能 [7]。

### 安全性考量與風險管理

#### 全面安全策略

PDF 資料提取 API 面臨多重安全挑戰，需要建立多層防護機制。檔案上傳功能容易受到各種安全威脅，包括惡意檔案注入、伺服器端請求偽造（SSRF）等 [8]。



Security Risk Assessment Matrix showing the relationship between risk level and implementation effort

關鍵安全措施包括：檔案類型驗證，防止惡意檔案傳輸；檔案大小限制，防止拒絕服務攻擊；惡意軟體掃描，檢測潛在威脅；輸入清理，確保檔案名稱不包含特殊字元或潛在安全威脅 [7]。

Adobe 的安全實踐提供了重要參考，其 PDF Extract API 採用深度防禦方法，包括使用者生成內容暫時快取、iframe 安全隔離、內容安全政策（CSP）等措施 [9]。

## 認證與授權機制

實施強固的認證系統對於保護 API 免受未授權存取至關重要。多層認證在雲端環境中特別有效，AES 演算法相較於 3DES 在執行時間、請求和回應時間方面表現更優，更適合實施多層認證<sup>[10]</sup>。

## 性能優化與擴展性設計

### 解析性能優化

PDF 解析性能直接影響 API 的使用者體驗。研究顯示，Grobid PDF 提取包配合 NLP 工具 Spacy 達到了 93% 的最高 f-1 分數，記憶體消耗最少，僅為 46.13 MB<sup>[11]</sup>。對於表格資料提取，pdfplumber 方法的識別率在高等教育品質報告中達到平均 96%，能有效處理空白和欄位錯位等複雜情況<sup>[12]</sup>。

### 快取與資源管理

實施快取機制對於提高 API 性能至關重要。可以快取經常存取的資源，如字體或圖像，消除冗餘的磁碟存取，從而加快後續 PDF 的生成時間<sup>[13]</sup>。對於資源密集型操作，如圖像處理和字體嵌入，應進行優化，包括壓縮圖像、使用適當的圖像格式、選擇性嵌入字體等<sup>[13]</sup>。

### 資料提取與處理策略

#### 多格式資料提取

現代 PDF 資料提取 API 需要處理多種資料類型，包括文字、表格、圖像和元資料。Adobe PDF Extract API 能夠提取文字區塊（段落、標題、清單、腳註等），包含字體、樣式和其他文字格式資訊<sup>[4]</sup>。表格會被提取和解析，包含內容和表格格式資訊，服務會自動識別跨多行或多列的表格儲存格<sup>[4]</sup>。

#### 結構化輸出格式

提取的內容應以結構化格式輸出，通常為 JSON 格式，表格可選擇性地輸出為 CSV 和 XLSX 檔案，圖像保存為 PNG 檔案<sup>[4]</sup>。這種結構化輸出使得資料能夠輕鬆存儲、分析和在各種下游系統中操作<sup>[4]</sup>。

### 特殊應用場景與解決方案

#### 多語言文件處理

對於多語言文件，特別是包含印尼語內容的社交媒體資料分析，可以利用 Twitter API 和 Pentaho Data Integration (PDI) 收集位置資料、推文資料和個人檔案資料<sup>[14]</sup>。使用 Streamlit 和 DuckDB 開發的互動式儀表板分析推文模式、主題標籤趨勢和食品飲料相關關鍵字，為企業提供有意義的洞察<sup>[14]</sup>。

#### 複雜表格查詢處理

對於包含複雜表格結構的 PDF 文件查詢，檢索增強生成（RAG）架構面臨挑戰。創新方法包括將 PDF 存儲在檢索資料庫中並單獨提取表格內容，對提取的表格進行上下文豐富化處理<sup>[15]</sup>。使用微調的 Llama-2-chat 語言模型在 RAG 架構內進行摘要，並通過 ChatGPT 3.5 API 使用單次提示增強表格資料的上下文意義<sup>[15]</sup>。

## 結論與建議

PDF 資料提取 Web API 的開發需要綜合考慮技術選擇、架構設計、安全防護和性能優化等多個面向。基於本報告的分析，建議採用混合解析策略，根據文件類型選擇最適合的解析工具，同時整合 AI 技術提升複雜文件的處理能力。在架構設計上，應遵循 RESTful 最佳實踐，實施完整的安全防護措施，並建立可擴展的快取和資源管理機制。未來發展方向包括進一步整合大型語言模型、增強 AI 驅動的分析能力，以及針對特定領域文件類型的專門化優化。這些建議將有助於開發人員和研究人員構建高效、安全、可靠的 PDF 資料提取服務，滿足現代數位化環境的多樣化需求。



1. <https://paperswithcode.com/paper/a-comparative-study-of-pdf-parsing-tools>
2. <https://www.theseattledataguy.com/challenges-you-will-face-when-parsing-pdfs-with-python-how-to-parse-pdfs-with-python/>
3. <https://onlinelibrary.wiley.com/doi/10.1002/jrsm.1732>
4. <https://developer.adobe.com/document-services/docs/overview/pdf-extract-api/>
5. <https://opensource.adobe.com/developers.adobe.com/apis/documentcloud/dcsdk/pdf-extract.html>
6. <https://daily.dev/blog/restful-api-design-best-practices-guide-2024>
7. <https://blog.poespas.me/posts/2024/05/28/handling-file-uploads-in-restful-apis/>
8. <https://tyk.io/blog/api-design-guidance-file-upload/>
9. [https://www.adobe.com/content/dam/cc/en/security/pdfs/AdobeDocumentServices\\_SecurityOverview.pdf](https://www.adobe.com/content/dam/cc/en/security/pdfs/AdobeDocumentServices_SecurityOverview.pdf)
10. <https://www.sciencepubco.com/index.php/ijet/article/view/24241>
11. <https://www.mdpi.com/2076-3417/12/3/1352>
12. <https://dl.acm.org/doi/10.1145/3696474.3696731>
13. <https://unidoc.io/post/optimizing-pdf-generation-high-volume-transactions/>
14. <https://ieeexplore.ieee.org/document/10761920/>
15. <https://www.semanticscholar.org/paper/4003b1bb28b3f8463accb3b5d57538708099f0ca>