# SAFEDOCS AI: PROTECTING SENSITIVE DATA IN DOCUMENT PROCESSING

Ensuring privacy and compliance through seamless data sanitization

By Victor Wu, Antoinette Davis

# THE CHALLENGE OF DATA PRIVACY IN DOCUMENT MANAGEMENT

- Handling sensitive data in documents is essential but challenging, especially with regulations like GDPR and CCPA.

- Organizations struggle to ensure data privacy while maintaining document usability.

- As organizations want to use more data in AI or other data intensive projects there are more risks like unintentional data leaks and regulatory fines.

- Manually redacting documents is inefficient and error-prone, especially at scale.

# THE SOLUTION: SAFEDOCS AI: A SECURE, SCALABLE SOLUTION

• SafeDocs AI automates sensitive data redaction using Azure AI's advanced OCR and language processing tools.

• Ensures document usability while scrubbing sensitive information.

•Provides value by keeping the data PII data free and intact to perform downstream tasks like AI model training, data analysis, and reporting.

• Allows compliance with GDPR, CCPA, and other privacy regulations.

# SAFEDOCS AI: CORE FUNCTIONALITIES

**Document Upload:** User-friendly upload interface supporting PDFs, **and images.**

**OCR and Text Extraction:** Azure AI **Computer Vision** extracts text and data.

**Sensitive Data Detection:** Uses Azure Language Service to detect PII like names, addresses, and credit card numbers.

**Data Scrubbing:** Automatically replaces sensitive information with placeholders.
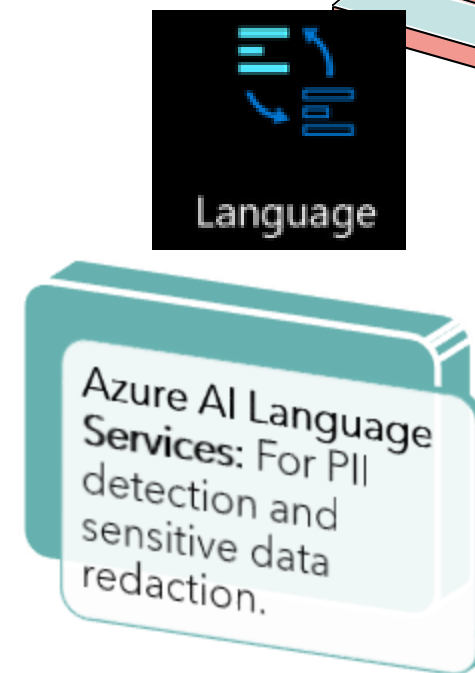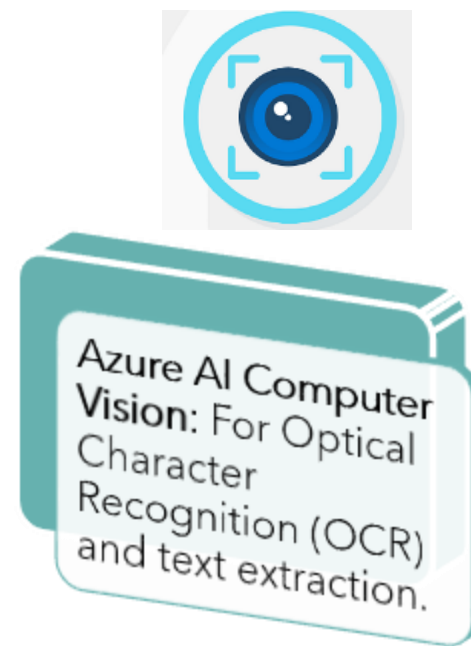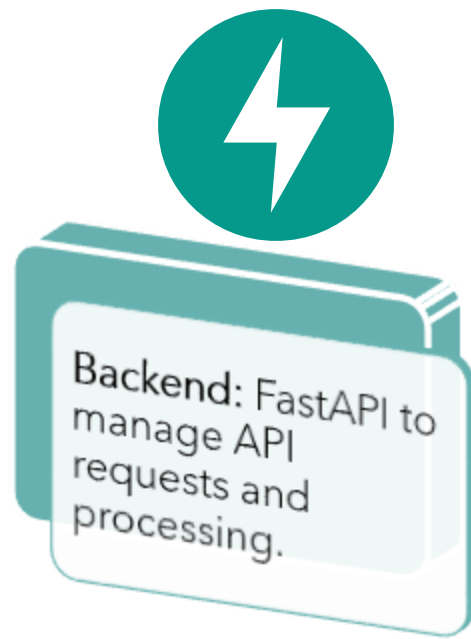
**Side-by-Side View:** Allows users to compare the original and sanitized documents.

**Data Summary:** Shows statistics on types and quantity of data scrubbed.

Frontend: React for an interactive and responsive user interface.

Backend: FastAPI to manage API requests and processing.

Azure AI Computer Vision: For Optical Character Recognition (OCR) and text extraction.

Language

Azure AI Language Services: For PII detection and sensitive data redaction.

# TECH STACK

# HOW WE THOUGHT THROUGH OUR APPROACH

**Sensitive Data Challenge:** Recognized the potential of data for AI, acknowledging that much of it includes sensitive information (PII).

**Compliance vs. Innovation:** Balanced the goal of responsibly leveraging data while ensuring GDPR and regulatory compliance.
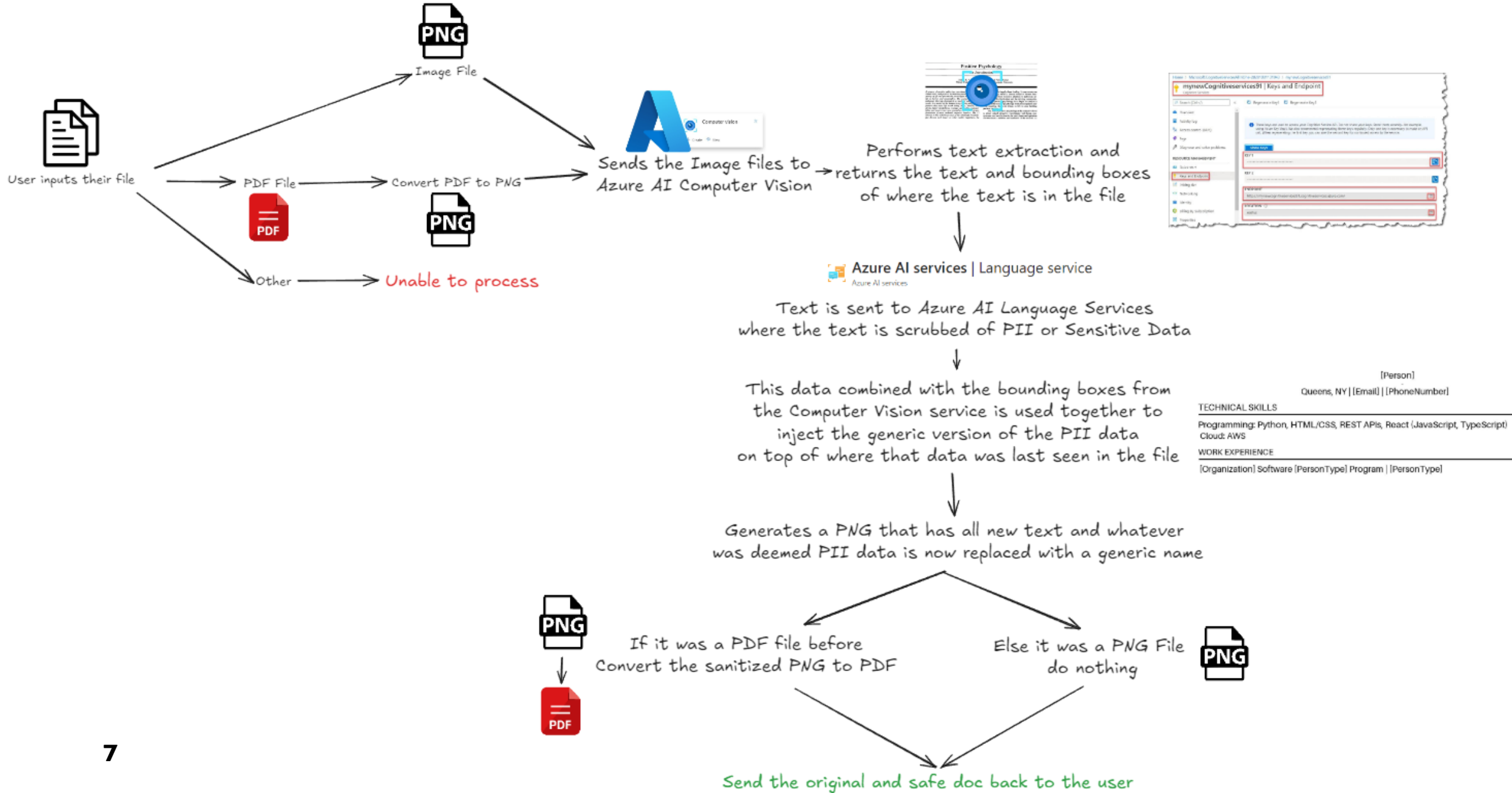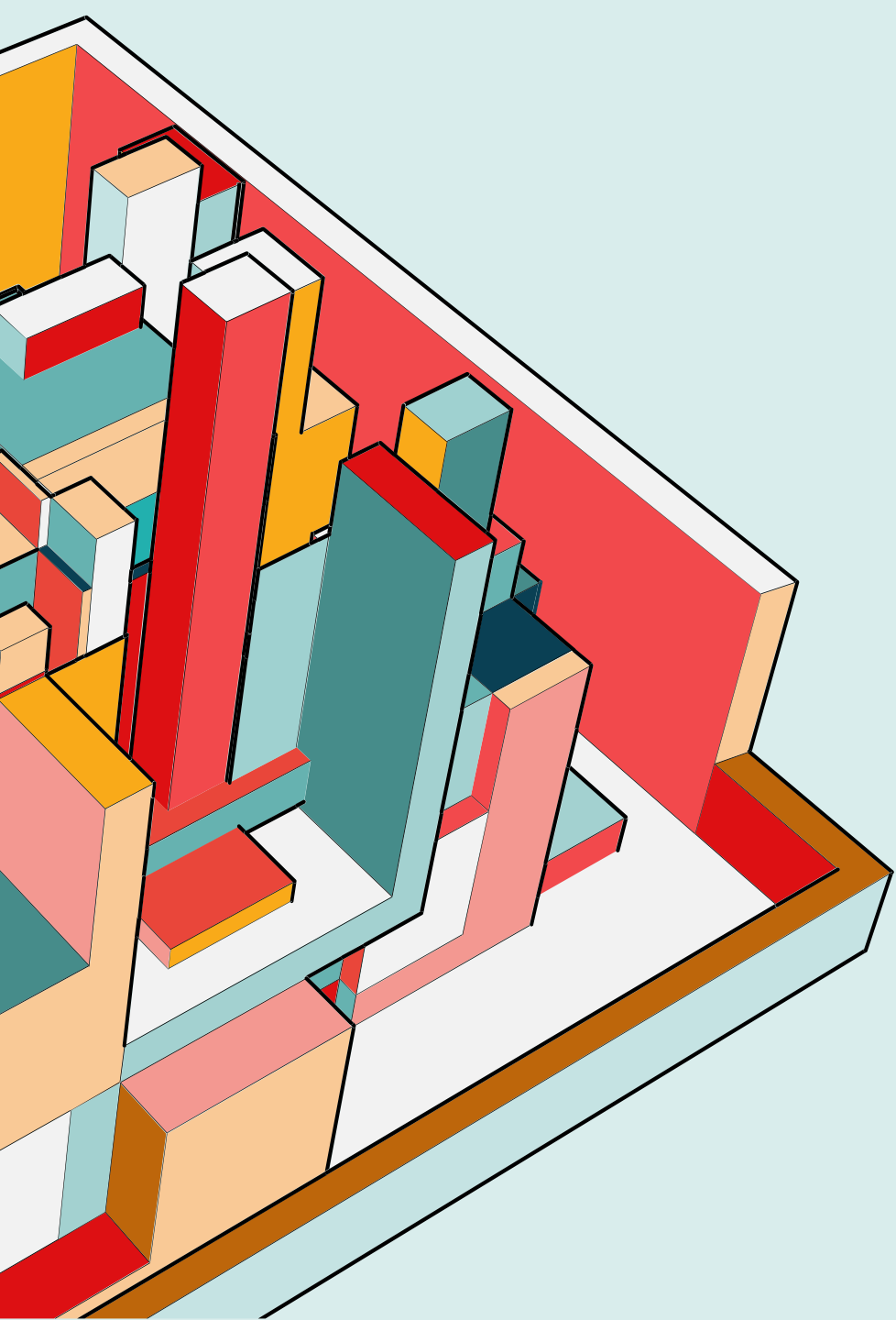
**Privacy-First Design:** Created a system prioritizing privacy with automated PII detection and redaction for compliant data use in training and analysis.

**Scalable Solution:** Built a scalable, efficient system for continuous data processing across sectors, securing sensitive information while preserving data utility.

# PROJECT ARCHITECTURE/WORKFLOW

PNG

Image File

User inputs their file

PDF File → Convert PDF to PNG

PNG

Other → Unable to process

Sends the Image files to Azure AI Computer Vision

Performs text extraction and returns the text and bounding boxes of where the text is in the file

Azure AI services | Language service
Azure AI services

Text is sent to Azure AI Language Services where the text is scrubbed of PII or Sensitive Data

This data combined with the bounding boxes from the Computer Vision service is used together to inject the generic version of the PII data on top of where that data was last seen in the file

Generates a PNG that has all new text and whatever was deemed PII data is now replaced with a generic name

PNG

If it was a PDF file before Convert the sanitized PNG to PDF

PDF

Else it was a PNG File do nothing

PNG

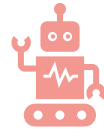Send the original and safe doc back to the user

LIVE DEMO

# IMPACT & USE CASES

**Business Impact:**
SafeDocs AI assists organizations in meeting GDPR and other regulatory standards, minimizing the risks of data breaches and associated fines.

**Document Processing for Compliance:** Ideal for legal, HR, and finance sectors where managing personally identifiable information (PII) is essential.

**AI Data Pipeline Safety:**
Pre-processes documents for privacy compliance, ensuring safe document handling for AI training workflows.

**AI Training Data Preparation:** Prepares documents for AI model training while safeguarding sensitive data.

**Document Analysis:**
Allows for in-depth document analysis without exposing PII.

**Regulatory Compliance:**
Simplifies redaction of sensitive information to meet various compliance needs.

**Client Documentation:**
Enables safe use of sanitized documents in client-facing workflows.

# CHALLENGES AND LEARNINGS

**Data Sensitivity:** Ensuring complete and accurate redaction of sensitive data, maintaining high precision to protect user privacy.

**Performance Optimization:** Addressing the challenge of balancing processing speed with redaction accuracy to handle large volumes of data without compromising quality.

**System Integration:** Integrating multiple Azure AI services into a unified backend system while managing complex API interactions and ensuring seamless communication between frontend and backend components.

**File Rebuilding:** Reconstructing sanitized files after data redaction to maintain original document structure and usability, a critical step for usability in downstream workflows.

**Frontend and Backend Integration Issues:** Resolving compatibility challenges between frontend and backend services to ensure smooth data flow and user experience.

# SAFEDOCS AI: FUTURE VISION AND CLOSING REMARKS

## Key Takeaways:

- SafeDocs AI enables secure document handling, ensuring compliance with privacy laws while preserving document utility
- Provides an efficient, scalable solution that integrates privacy into data pipelines.

## Future Enhancements:

- **Expanded File Support**: Include more file types (e.g., DOCX, XML).
- **Enhanced Redaction Accuracy**: Improve sensitivity detection and placeholder customization.
- **Analytics Dashboard**: Visualize trends in data scrubbing for deeper insights.
- **User Authentication**: Add security layers for authorized access.
- **Seamless AI Integration**: Enable direct AI processing on sanitized documents.