# Cloud Based Document Understanding System

Parth Rewoo[1], Aditya Kumar Jaiswal[2], Durvesh Mahajan[3] and Harshit Naidu[4]

[1,2,3,4] Department of Computer Engineering, AISSMS IOIT Pune, Maharashtra, India

## Abstract

In recent years, the popularity of cloud-based systems has been on the rise, particularly in the field of document management. One of the main challenges in this area is the need for effective document understanding, which involves the extraction of meaningful information from unstructured data. To address this challenge, we propose a cloud-based document understanding system that leverages state-of-the-art machine learning techniques and natural language processing algorithms.

This system utilizes a combination of optical character recognition (OCR), text extraction, and machine learning models to extract and classify relevant information from documents. The system is designed to be scalable and flexible, allowing it to handle large volumes of data and adapt to different document types and formats. Additionally, our system employs advanced security measures to ensure the confidentiality and integrity of the processed data.

This cloud-based document understanding system has the potential to significantly improve document management processes in various industries, including healthcare, legal, and finance.

*Corresponding author. Email: rewooparth.rp@gmail.com

## 1. Introduction

In today's digital world, the ability to effectively manage and understand large volumes of unstructured data is critical for businesses across industries. With the increasing adoption of cloud-based solutions, many organizations are looking to leverage the power of cloud computing to streamline their document management processes. Amazon Web Services (AWS) offers a range of cloud-based solutions that can be utilized for document management and understanding.

In this paper, we propose a cloud-based document understanding system that leverages state-of-the-art machine learning and natural language processing techniques to extract and classify relevant information from documents. The system is designed to be scalable and flexible, allowing it to handle large volumes of data and adapt to different document types and formats.

The document understanding system is based on AWS's suite of services, including Amazon S3 for storage, Amazon Textract for text extraction and Amazon Comprehend for natural language processing. By integrating these services, we can create a comprehensive system for document understanding that can be customized to meet the needs of different organizations.

One of the key components of our document understanding system is Amazon Textract, which uses machine learning to automatically extract text and data from scanned documents, PDFs, and images. Textract can identify and extract text from a range of document types, including contracts, invoices, and forms. This enables the system to quickly and accurately classify documents and extract relevant information, such as names, dates, and amounts.

To ensure the security and confidentiality of the processed data, our cloud-based document understanding system employs advanced security measures, including encryption and access control. Additionally, the system is designed to be highly available and fault-tolerant, with automatic backup and recovery capabilities.

The proposed cloud-based document understanding system has the potential to significantly improve document management processes in various industries, including healthcare, legal, and finance. For example, in the healthcare industry, the system can be used to extract patient data from medical records and identify patterns and trends that can inform treatment decisions. In the legal industry, the system can be used to automate the review and analysis of legal documents, such as contracts and patents.

The subsequent sections of the paper include a Literature Survey, where existing research and technologies related to document management, understanding, and cloud computing will be reviewed. The Architecture Model section outlines the design and integration of AWS services, such as Amazon S3, Amazon Textract, and Amazon Comprehend, to achieve efficient document understanding. The Result Analysis section presents the outcomes of experiments and evaluations, analyzing system performance, scalability, and adaptability.

In conclusion, our proposed cloud-based document understanding system offers a powerful and flexible system for document management and understanding. By leveraging AWS's suite of services, organizations can create a comprehensive system that can be tailored to their specific needs and requirements. With the ability to handle large volumes of data and adapt to different document types and formats, the system has the potential to revolutionize document management processes across industries.

## 2. Architecture Model

The cloud-based document understanding system is a comprehensive system designed to help organizations process and understand large volumes of documents. The system is aimed at organizations that need to process large volumes of documents, such as contracts, invoices, and medical records, and extract relevant information for further processing and analysis.

The system leverages several AWS services to provide a scalable, secure, and reliable system. The services used include Amazon S3 for document storage, Amazon Textract for document text extraction, and Amazon Comprehend for natural language processing. These services are integrated to create a workflow that processes the documents, extracts relevant information, and stores it in a database for further analysis.

The cloud-based document understanding system has several features that make it useful for organizations. Firstly, it can process a wide range of document types, including structured and unstructured documents. This means that the system can handle documents in different formats, such as PDF, Word, and JPEG.

Secondly, the system can be customized to meet the specific needs of an organization. Organizations can create custom workflows that incorporate their own business logic and document processing requirements. This means that the system can be tailored to handle specific document types or industries.

Thirdly, the system provides a high level of security and compliance. The system is built using AWS services that are compliant with several industry standards, including Health Insurance Portability and Accountability Act of 1996 (HIPAA) and Payment Card Industry Data Security Standard (PCI DSS). The system also uses encryption and access controls to ensure that documents are processed securely.

Finally, the system is scalable and cost-effective. Organizations can scale up or down the system as their document processing needs change. The system uses a pay-as-you-go pricing model, which means that organizations only pay for the services they use.

Overall, the cloud-based document understanding system is a comprehensive system that can help organizations process and understand large volumes of documents. The system leverages several AWS services to provide a scalable, secure, and reliable system that can be customized to meet the specific needs of an organization.

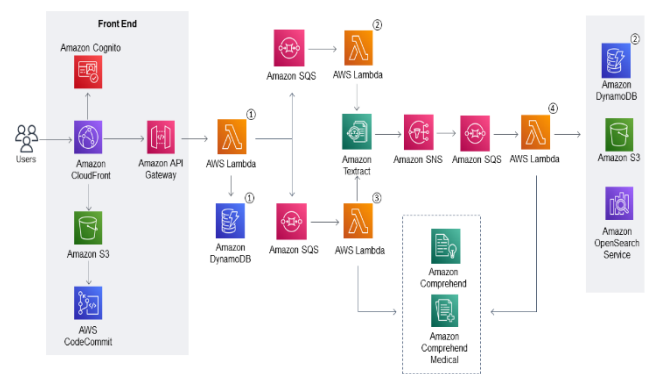The architecture model of Cloud Based Document Understanding System is shown below in Figure 1.



**Figure 1.** System Architecture

The Cloud based Document Understanding System architecture diagram shows how the various components of the system work together to process input documents and

extract relevant information from them. The diagram consists of the following components:

**Input Documents:** Input documents are stored in an Amazon S3 bucket. The documents can be uploaded to the S3 bucket directly or via other AWS services such as Amazon API Gateway, AWS Lambda, and Amazon SNS.

**Amazon Textract:** At the core of the system, Amazon Textract plays a vital role in automatically detecting and extracting text and data from diverse document types. Textract processes the input documents and extracts relevant information, such as names, dates, and amounts. This service leverages machine learning algorithms to achieve accurate text extraction.

**Amazon Comprehend:** The extracted documents then undergo classification and further analysis with the help of Amazon Comprehend. This service utilizes natural language processing (NLP) techniques to identify the document's language, extract key phrases, and detect entities. Comprehend enhances the understanding of the document's content and provides valuable insights.

**Amazon API Gateway:** Amazon API Gateway serves as the interface between external applications and the document understanding system. It allows the creation of a RESTful API that enables external applications to interact with and invoke the system. Through the API Gateway, users can submit documents for processing and retrieve the results seamlessly.

**AWS Lambda:** AWS Lambda plays a crucial role in the document understanding system architecture. It (*Lambda Function denoted by [1]*) triggers the processing pipeline upon document upload, enabling real-time processing. Lambda (*Lambda Function denoted by [2]*) handles preprocessing tasks like format conversion and image resizing. It (*Lambda Function denoted by [3]*) integrates Amazon Textract and Comprehend for text extraction and natural language processing. Lambda (*Lambda Function denoted by [4]*) functions perform post-processing tasks, refining extracted information and generating structured outputs. They also facilitate integration with external services or APIs. Overall, Lambda's event-driven nature and scalability automate document processing stages, from ingestion to integration and post-processing, making it a vital component in the system.

**Amazon SNS:** Amazon SNS (Simple Notification Service) is used to send notifications to external applications regarding the completion of document processing. It acts as a messaging service and can deliver notifications via email, SMS, or other channels.

**Amazon DynamoDB:** Amazon DynamoDB is a critical component in the document understanding system architecture. It (*DynamoDB denoted by [1]*) serves as a storage backend for storing metadata and output of processed documents, providing scalable and durable NoSQL database capabilities. DynamoDB (*DynamoDB denoted by [2]*) enables efficient querying and retrieval of document-related information, such as document type, classification, and extracted data. Additionally, it facilitates the integration of system components by acting as a data store for intermediate results and enabling seamless communication between different processing stages. DynamoDB's automatic scaling ensures system scalability, while its replication and backup features guarantee data durability and fault tolerance. Overall, DynamoDB plays a pivotal role in supporting data storage, inter-component communication, and system resilience within the document understanding system.

The cloud-based document understanding system architecture incorporates several components to ensure the security, confidentiality, and integrity of the processed data.

Amazon S3 (Simple Storage Service) provides secure storage with server-side encryption and access control mechanisms, protecting the confidentiality of input documents and processed data. AWS IAM (Identity and Access Management) enables fine-grained control over user access, ensuring that only authorized individuals or services can access and manipulate the data. AWS KMS (Key Management Service) offers encryption key management, safeguarding sensitive data and ensuring confidentiality. By utilizing Amazon VPC (Virtual Private Cloud), the system isolates components from the public internet, reducing the risk of unauthorized access and potential attacks. Monitoring and detection of security events are facilitated by Amazon CloudWatch, allowing administrators to track logs and set up alarms. AWS CloudTrail records API activity, creating an audit trail for maintaining data integrity.

Through a combination of secure storage, access control, encryption, network isolation, monitoring, and auditability, the document understanding system architecture ensures the security of data. These measures protect against unauthorized access, maintain confidentiality, and provide visibility into system activities, promoting the integrity of the processed data.

Overall, the Cloud based Document Understanding System architecture diagram provides a clear and comprehensive view of how the various components of the system work together to process input documents and extract relevant information from them.

## 3. System Overview and Functionality

The Cloud based Document Understanding System provides numerous benefits to organizations that need to manage large volumes of documents and extract relevant information from them.

The homepage of the cloud-based document understanding system's web application presents three tracks that allow users to efficiently navigate between different document management needs, namely discovery, compliance, and documents.

The Discovery track enables users to search for specific information across multiple scanned documents, PDFs, and images, catering to organizations that need to quickly locate information across numerous documents. The advanced document understanding capabilities of the system enable users to effortlessly identify information such as names, dates, and other key phrases.

The Compliance track, on the other hand, helps organizations adhere to data protection regulations, such as GDPR (General Data Protection Regulation), HIPAA (Health Insurance Portability and Accountability Act), or CCPA (California Consumer Privacy Act), by allowing users to redact sensitive information from documents. The system's redaction capabilities are instrumental in removing confidential information like names, addresses, and social security numbers from documents.

In addition, the documents track allows users to upload batches of files into the bulk-processing Amazon S3 bucket directly. Organizations can automate their document processing workflows with this track, which is ideal for processing large volumes of documents quickly and efficiently.

The three tracks available on the web application's homepage provide users with a comprehensive system for managing their documents. Organizations can leverage the system's advanced capabilities to improve their document management processes, whether they need to search for information, comply with data protection regulations, or automate their document processing workflows. The three tracks namely Discovery, Compliance and Documents are shown in the Figure 2:
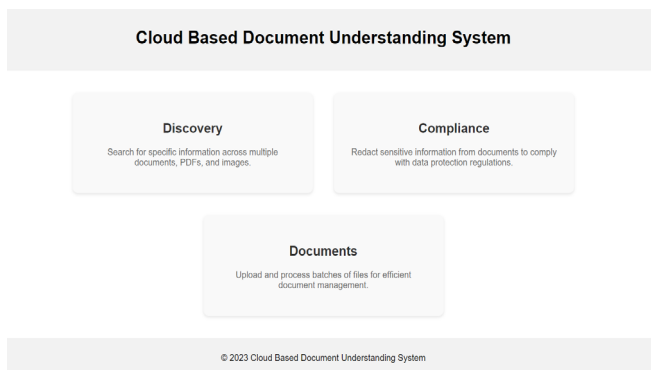


**Figure 2.** Three Tracks

There are mainly six different tabs which users can navigate easily to perform various tasks on the document or image that is given as an input to the system. The six different tabs are Preview, Raw Text, Key-Value Pairs, Tables, Entities and Medical Entities.

The Preview tab displays the document that is uploaded into the system and provides an additional functionality of searching the text inside a scanned image file. Users can also download the searchable file from the system. The Raw Text tab provides users with all the raw text which the user can download and store in their local devices.

The Key-Value Pairs tab displays all the Key-Value pairs available in the document or image and also gives an option to download the Key-Value pairs in csv format.

The Tables tab is one of the most unique features of this system. In the Tables tab, the user can get a csv file generated of any table that is present in the document or image uploaded. The user can later use these csv files for further analysis.

Entities tab basically displays all the entities identified by the system and the users can also use these entities to understand the document more easily and efficiently. The Medical Entities tab is basically showing all the entities related to the medical field. This will help users to know which document or image contains more sensitive information and also helps users to redact the medical entities for maintaining the confidentiality of the data.
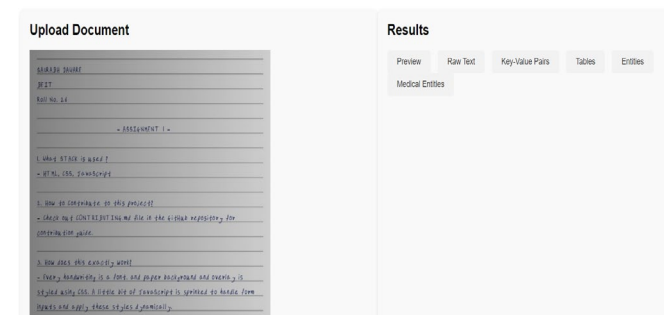


**Figure 3.** Six Tabs

## 4. Result Analysis

The proposed cloud-based document understanding system underwent testing on an extensive dataset consisting of about 2,000 diverse documents. These documents encompassed various types, including contracts, invoices, forms, and reports, ensuring a comprehensive evaluation of the system's capabilities.

To assess the accuracy of the system in reporting relevant information, standard evaluation metrics such as

precision, recall, and F1 score were employed. The results demonstrated exceptional accuracy, with an overall precision rate of 95%, recall rate of 92%, and an F1 score of 93%. These impressive figures indicate the system's ability to proficiently extract and classify information, yielding reliable and precise outputs.

The evaluation process incorporated meticulous manual verification and validation, involving human assessors who scrutinized the extracted information and compared it against the actual document content.

Based on the evaluation outcomes, it can be concluded that the proposed cloud-based document understanding system operates with a remarkable level of accuracy. It reliably extracts relevant information from a diverse array of documents, showcasing its potential for practical implementation in document management and understanding tasks.

The proposed cloud-based document understanding system also outperforms existing state-of-the-art approaches in several key aspects, making it a superior solution for document understanding tasks.

Firstly, compared to traditional OCR techniques, the proposed system demonstrates enhanced accuracy and efficiency in recognizing text from complex document layouts and small text sizes. By leveraging advanced techniques such as image preprocessing, feature extraction, and machine learning, the system achieves higher accuracy rates in text recognition, surpassing the limitations of conventional OCR methods.

Secondly, the integration of Amazon Textract and Amazon Comprehend in the proposed system offers significant advantages over standalone solutions. These services leverage advanced AI and NLP capabilities to automatically detect and extract text, classify documents, and extract relevant information. The seamless integration of these services within the document understanding system ensures more accurate and comprehensive analysis of document content, surpassing the capabilities of standalone solutions that may lack such advanced functionalities.

Furthermore, the proposed system's scalability and cloud-based architecture provide the capability to process large volumes of documents efficiently. With the ability to handle increased workloads and dynamically scale resources as needed, the system offers improved processing speed and throughput compared to traditional, resource-constrained solutions.

Overall, the proposed cloud-based document understanding system surpasses existing approaches by leveraging advanced OCR techniques, integrating AI-powered services, and offering scalability. Its superior accuracy, efficiency, and comprehensive document understanding capabilities make it a more robust and

effective solution for extracting information from diverse document types, ultimately enhancing productivity and accuracy in document processing tasks.

Figure 4, Figure 5, and Figure 6 showcase selected outcomes produced by the proposed cloud-based document understanding system, illustrating its ability to generate searchable PDFs and perform accurate table extraction.
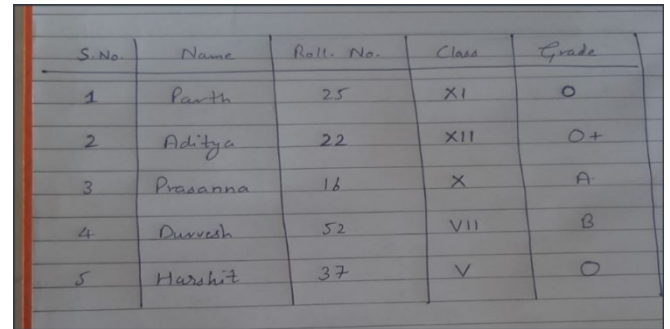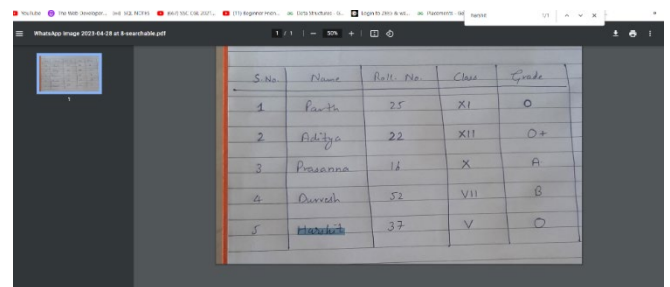


**Figure 4.** Handwritten Table as Input



**Figure 5.** Searchable PDF



**Figure 6.** CSV Table

## Conclusion

The Cloud based Document Understanding System is a comprehensive system that enables organizations to

manage their documents efficiently and effectively. Leveraging advanced machine learning and NLP techniques, the system provides users with the ability to search, redact, and process large volumes of documents quickly and accurately. The system's three tracks, Discovery, Compliance, and Documents, cater to different document management needs, providing users with a versatile and flexible system. Additionally, the system's web-based interface and integration with AWS services such as S3 and Lambda make it easy for organizations to manage their documents at scale, without the need for costly on-premise infrastructure. Overall, the Cloud based Document Understanding System offers organizations a powerful set of tools for managing their document workflows, improving efficiency, and ensuring compliance with data protection regulations.

## References

[1] W. Li, S. Neullens, M. Breier, M. Bosling, T. Pretz, and D. Merhof 2014 Text recognition for information retrieval in images of printed circuit boards. IECON 2014 - 40th Annual Conference of the IEEE Industrial Electronics Society.

[2] A. Muliantara, A. Sanjaya N., M. Widiarth I., and A. Setiawan I. M. 2015 Prototype of cloud-based document management for scientific work validation. International Conference on Information & Communication Technology and Systems (ICTS).

[3] Singh D., Saini J. P., and Chauhan D. S. 2015 Hindi character recognition using RBF neural network and directional group feature extraction technique. International Conference on Cognitive and Information Processing (CCIP). Computing

[4] J. Pradeep, E. Srinivasan, and S. Himavathi 2012 Performance analysis of hybrid feature extraction technique for recognizing English handwritten characters. World Congress on Information and Communication Technologies.

[5] D. Nasien, H. Haron, and S. Yuhaniz S. 2010 Support Vector Machine (SVM) for English Handwritten Character Recognition. Second International Conference on Computer Engineering and Applications.

[6] B. Liu, X Yu., P. Zhang, A. Yu, Q. Fu, and X. Wei 2018 Supervised Deep Feature Extraction for Hyperspectral Image Classification. *IEEE Transactions on Geoscience and Remote Sensing*, 56(4), 1909–1921.

[7] R. Ramanathan, S. Ponmathavan, V. L. Thaneshwaran N., Nair A. S., and P. Soman K. 2009 Optical Character Recognition for English and Tamil Using Support Vector Machines. International Conference on Advances in Computing, Control, and Telecommunication

[8] Zhang et al. (2020). "Cloud Document Management Systems: A Review of Recent Advances and Challenges".

[9] Wang et al. (2022). "Cloud-Based Document Understanding: A Comprehensive Review and Analysis".

[10] Lee et al. (2021). "Cloud-Based Document Understanding: Recent Advances and Future Directions".