

# A Study on Intelligent Document Processing Using AWS

**Smt. K.S. Sukrutha<sup>1</sup>, Ms. Harini.S<sup>2</sup>, Ms.Kusuma. M. V.<sup>3</sup>**

<sup>1</sup>Assistant Professor, Department of Computer Science, M.M. K & S.D.M Mahila Maha Vidyalaya,  
Mysuru, India

<sup>2,3</sup>III BCA Students, M.M. K & S.D.M Mahila Maha Vidyalaya, Mysuru, India

## ABSTRACT

This article mainly approaches the topic of Intelligent Document Processing from the viewpoint of the users of cloud computing platforms and the end-users. The market for commercial OCR, document categorization and data extraction technologies are briefly reviewed to the extent that it is publicly available. The necessity for an effective and efficient retrieval of the stored information is increased by digital repositories. In this study, we suggest that the phases of document layout analysis, document picture categorization and understanding on digital documents should intensively apply intelligent techniques. Specifically, the Intelligent Document Processing can be used instead of Artificial Intelligence, Deep Learning and RPA to retrieve the data and convert to understandable form.

## KEYWORDS

Intelligent Document Processing (IDP), Amazon Web Services (AWS), Robotic Process Automation (RPA), Straight Through Processing (STP), Natural Language Processing (NLP), Personally Identifiable Information (PII), Deep Learning(DL).

## INTRODUCTION

Intelligent document processing (IDP) is the process of converting the unstructured data like email, images and PDF document to usable data. Usually all 80%-90% of data in the business company is in the form of unstructured format. Intelligent document processing is one of the processes that is carried out by artificial intelligence and AI technologies and one of the best automations in coming up generation it involves computer vision, Deep Learning and machine learning, etc. to extract data [2]. Intelligent document processing crosses 175 zettabytes all over the world by 2025 with the document from PDF 's, email etc.

In order to convert the unstructured and semi structured data to usable information IDP mainly act as a key to the RPA. Without intelligent document processing RPA-automation process requires a knowledge worker to read documents and extract data. Further, if RPA alone is used for the conversion of structured, unstructured and semi structured data it leads to less productivity, less economic benefits, less accurate and less customer satisfaction [5].

Together RPA and IDP will provide a very useful tool to automate a business company or enterprises. RPA and intelligent document processing which are very much friendly to user and non-invasive are widely applicable across all companies. [2]

Intelligent document processing provides

1. Save cost- reduces much processors to process large volume of data.
2. Straight through processing (STP) - minimize the need of knowledge worker when intelligent document processing is used.
3. Easy to use - any business or enterprise can easily use automate the process.
4. It provides efficiently – it provides document centric process.
5. High accuracy – when AI is used obviously the accuracy is increased.
6. Strategic goals boosted-when a business uses the intelligent document processing automatically their goals like improving customer experience is also increased.
7. Enterprise with goals like superior customer experience and to provide low cost are opting them to intelligent document processing.
8. Though the work of intelligent document processing is little bit similar as OCR and RPA many of them get confused with all of them but they all perform differently.
9. Intelligent document processing does more than OCR in which OCR converts the bitonal imagery to a machine-readable form (bitonal imagery means single bit images that has only 0 or 1 (0 to 255) i.e., white or black but the intelligent document processing process the unstructured data to a structured one) [2]
10. Unlike OCR, intelligent document processing uses Artificial Intelligence (AI) and machine learning (ML) technologies to change or transform unstructured and semi-structured material into data that can be understood. Intelligent document processing primarily uses robotic process automation (RPA) to extract data, improve validation, and automatically enter information into current applications[4].
  - a. **Machine Learning** - Machine learning emphasizes primarily on the study of computer algorithms that are enhanced by data-driven experiences.
  - b. **Artificial Intelligence** – Unlike natural intelligence, which is used by living things, artificial intelligence is a type of intelligence that is implemented on computers.
  - c. **Deep Learning** - Deep Learning is a combination of many machine learning techniques that primarily makes use of multiple layers in diverse neural network architectures.

## INTELLIGENT DOCUMENT PROCESSING IN AWS

Companies like Amazon uses intelligent document processing (Amazon A2I) to process some critical document and for Natural Language Processing with Amazon comprehend and Amazon A2I.



**Figure 1. Stages of IDP workflow**

Figure 1 mainly explains various stages of the intelligent document processing pipeline and the connectivity between each steps starting from the application or document submission to investigation and closing the application or documents. We can also observe the technical details such as data capture classification and extraction stages and document enrichment, review and verification and extend the solution to provide analytics and visualization for a claim’s fraud use case.

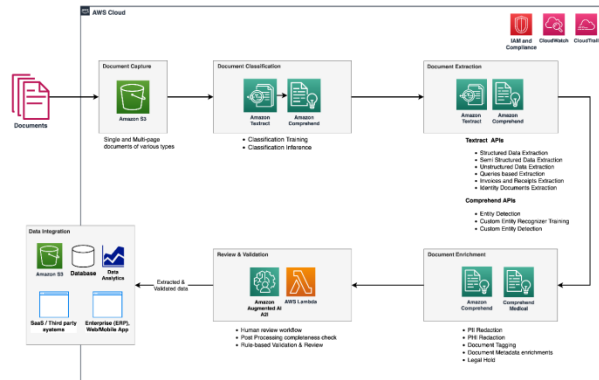


Figure 2. Architecture of IDP

The above architecture diagram explains the various stages of IDP workflow

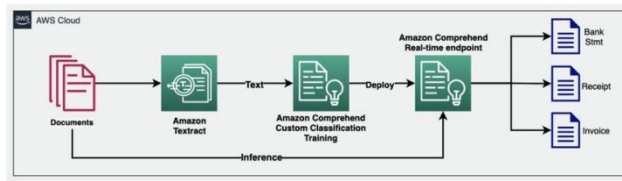
- DATA CAPTURE**- it usually centralizes the data and claims.
- CLASSIFICATION**- it usually sends all unstructured and semi-structured data through the pipeline to extract the data.
- EXTRACTION**-it extracts all the information from the claim and tag the notes so that search will be easy.
- ENRICHMENT**-all the unstructured and semi- structured are all identified and made a proper view to the readers.
- REVIEW AND VALIDATION**-for the review and the validation the converted document is sent to the authenticated person to check it.
- READY TO USE**-the converted file is updated to the main database within 24 hours after getting from the validation from authenticated person.

## PHASES IN INTELLIGENT DOCUMENT PROCESSING

- CLASSIFICATION PHASE** – Collected documents of various types are categorized before further extraction through Amazon Comprehend custom classification which is a two step process as shown illustrated in figure 3. The custom classification is the process mainly helps to automate the document and identify the missing document from the packet.

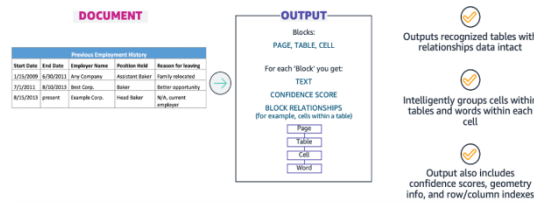
## TWO CUSTOM CLASSIFICATION STEPS

- EXTRACT TEXT USING AMAZON TEXT EXTRACT** – From any document or image, Amazon textextract uses machine learning to automatically extract text, handwriting, image and data. Instead of hours or days, Textextract can extract the data in minutes. Thus, we can process the documents quickly and extract the information out of it.
- TRAIN AMAZON COMPREHEND CUSTOM** - training Amazon comprehend custom classification or documents classifier to recognize the classes of interest based on the text content [1].



**Figure 3. Document classification**

2. **EXTRACTION PHASE** - in the extraction phase, we mainly extract the document where we have the data using Amazon comprehend
  - We also use the claims processing packets to process it.
  - CMS-1500 CLAIM FORM is used to extract data.
  - Non-institutional provide a CMS-1500 from as a standard claim.
  - CMS-1500 should be processed accurately or else it can slow down the claim process or delay payments by the carrier.[1]



**Figure 4. Data extraction process**

**THE KEY SERVICES OF INTELLIGENT DOCUMENT PROCESSING**

- 1) **AMAZON Textract** - it mainly extracts the hand writing, data from scanned document.
  - It is beyond the OCR (Optical Character Recognition)
  - Used mainly to extract data from forms and tables and understand it.
  - It includes machine learning to read and process the document, especially the non-manual form data or document.[3]
- 2) **AMAZON COMPREHEND**-mainly identifies the quantity, location, person date, dominant language, Personally Identifiable Information (PII)
  - It classifies intro-relevant classes.
  - It is a natural language processing (NLP) service that use Machine Learning(ML) to extract the insights from text[3].
- 3) **AMAZON AUGMENTED AI (AMAZON A21)** - it mainly provides a bridge between Amazon textract and Amazon comprehend to provide a ability to introduce a human review or validation within intelligent document processing flow. It mainly uses the machine learning service to make easy to build workflow required for human review [3].

## CONCLUSION

In this study, we saw how the structured, unstructured and semi-structured data, files, claims processed in AWS AI services and to automate the intelligent document processing pipeline. Here, we emphasized an idea to classify the documents into various document classes using an Amazon comprehend custom classifier, and to use Amazon to extract unstructured, semi-structured, structured and specialized document types.

The classification and extraction phase are expanded with Amazon textract. We also use Amazon comprehend pre-defined entities and custom entities to enrich the data and show how to extend the intelligent document processing pipeline to integrate with analytics and visualization services for further processing.

## REFERENCES

1. Intelligent document processing with AWS AI services: Part1|AWS Machine Learning Blog, <https://aws.amazon.com/blogs/machine-learning/part-1-intelligent-document-processing-with-aws-ai-services/>
2. Content from a website, <https://www.automationanywhere.com>
3. Intelligent document processing AWS Solutions for Machine learning AI/ML, <https://aws.amazon.com>
4. Intelligent Document Processing-Methods and Tools in the real World(Published paper), <https://www.researchgate.net> - Graham A Cutting, Independent Researcher, F [grahamcutting@cantab.net](mailto:grahamcutting@cantab.net); Anne-Francoise Cutting - Decelle : Universite de Geneve / CUI, CH ; [anne-francoise.cutting-decelle@unige.ch](mailto:anne-francoise.cutting-decelle@unige.ch)
5. Deloitte, <https://www2.deloitte.com>