



**PROJECT REPORT**

**OF**

**PASSION PROJECT**

**“INTELLIGENT DOCUMENT PROCESSING (IDP)”**

**UNDERTAKEN AT**

**SABUDH FOUNDATION, MOHALI**

**SUBMITTED IN FULFILLMENT OF THE**

**DATA SCIENCE INTERNSHIP**

**Interns:**

**Aditya Thite, Ali Haider Khan, Ishleen Kaur, Krishnendu Biswas**

**Mentored by:**

**Dr. Chandan Biswas, Mr. Manpreet Singh, Mr. Akshay Kumar, Mr. Gopi Maguluri**

---

**SABUDH FOUNDATION MOHALI**

PhaseVIII A, Plot No. C, 184, Industrial Area, Sector 75,  
Sahibzada Ajit Singh Nagar, Punjab 160071

Email: [admin.info@sabudh.org](mailto:admin.info@sabudh.org), Mob: +91 8837662054

## ABOUT THE ORGANIZATION



Sabudh Foundation is formed by the leading data scientists in the industry in association with the Punjab government with the objective to bring together data and young data scientists to work on focused, collaborative projects for social benefit. We aim to enable the youth to use powerful AI technologies for the greater good of society by working on real-world problems in partnership with nonprofit organizations and government agencies, to tackle data-intensive high impact problems in education, healthcare, public policy, agriculture etc.

## ABSTRACT

Intelligent Document Processing (IDP) is revolutionizing the way businesses handle data extraction and document management. By automating the process of manual data entry from paper-based documents or document images, IDP integrates seamlessly with other digital business processes, offering a range of benefits including scalability, cost-efficiency, and improved customer satisfaction. IDP solutions allow businesses to accurately scan and process large volumes of documents, significantly reducing human errors and enhancing efficiency, ensuring that companies can manage heavy operational demands with improved accuracy and efficiency. Additionally, automating document processing and analysis reduces overhead costs associated with manual data entry, eliminating repetitive tasks and overcoming bottlenecks, thereby boosting productivity and streamlining workflows. Furthermore, IDP accelerates the handling of customer documents, automating tasks such as customer onboarding, bookings, and payments, leading to enhanced processing speed and accuracy, which in turn improves customer service and relationships. Key subtasks in IDP include Table Structure Recognition (TSR) for processing structured data from tables, Named Entity Recognition (NER) for extracting specific data points from unstructured text, and Text Summarization (TS) for generating concise summaries from longer documents. IDP has diverse applications across multiple industries: in healthcare, it improves the management of healthcare records by extracting data from patient documents and organizing medical records; in finance, it automates expense management, invoice processing, and streamlines payment processes; in the legal sector, it aids in analyzing contracts and extracting critical data from legal documents; in logistics, it processes documents related to shipments and transit permits, reducing human errors and speeding up logistic functions; and in human resources, it extracts information from resumes, manages payroll, and handles leave allotments, allowing HR teams to focus on selecting top candidates and managing HR functions more efficiently. The implementation of IDP solutions in various business sectors demonstrates the transformative potential of automating document processing, leveraging machine learning and artificial intelligence to enhance operational efficiency, reduce costs, and improve customer satisfaction. This project aims to explore and develop advanced IDP solutions tailored to the specific needs of different industries, driving digital transformation and operational excellence.

## ACKNOWLEDGEMENT

We are writing to express our deep gratitude towards the esteemed members of Sabudh Foundation, Mohali, and Tatra Data whose unwavering support and dedication made the training phase successful and lucrative for us and our fellow interns. Learning from experienced professionals like **Dr. Sarabjot Singh Anand** (Co-Founder and Head of Academics, Sabudh Foundation & Co-Founder, Tatra Data), **Mr. Partha Sarthi Mukherjee**, **Mr. Prabhnoor Singh**, **Ms. Suparna Dutta**, **Dr. Shafila Bansal** and getting guidance and mentorship from **Dr. Chandan Biswas**, **Mr. Akshay Kumar**, **Mr. Manpreet Singh** during different phases of Internship and project development was all an honour.

# CONTENTS

<b>1.</b>	<b>INTRODUCTION (IDP)</b>	<b>(6-7)</b>
1.1.	ABOUT IDP	
1.2.	SUBTASKS	
1.3.	APPLICATIONS	
<b>2.</b>	<b>RESEARCH UNDERWENT</b>	<b>(8-15)</b>
2.1.	LITERATURE REVIEW	
2.1.1.	NER	
2.1.2.	TSR	
2.1.3.	TS	
<b>3.</b>	<b>PROJECT DEVELOPMENT</b>	<b>(16-52)</b>
3.1.	NER TASK	
3.1.1.	MODEL USED .	
3.1.2.	RESULTS/COMPARISON WITH OTHER MODELS	
3.2.	TSR TASK	
3.2.1.	MODEL SELECTION	
3.2.2.	RESULTS/COMPARISON	
3.3.	TS TASK	
3.3.1.	MODEL SELECTION	
3.3.2.	RESULTS/COMPARISON	
<b>4.</b>	<b>USER INTERFACE</b>	<b>(53-55)</b>
<b>5.</b>	<b>CONCLUSION AND FUTURE SCOPE</b>	<b>(56-57)</b>
	<b>CITATIONS AND REFERENCES</b>	

# CHAPTER 1: INTRODUCTION (IDP)

## 1.1 ABOUT IDP

Intelligent document processing (IDP) is automating the process of manual data entry from paper-based documents or document images to integrate with other digital business processes. DP offers a range of benefits for businesses. The following are some of the key advantages.

- **Scalability:** Manual document processing can result in human errors, reducing the efficiency of your business. It also limits how many documents you can process at a time. With IDP solutions, you can accurately scan documents at scale. ML/AI solutions process documents without mistakes. You can manage heavy operational demands with improved accuracy and efficiency.

- **Cost-efficiency:** Automation of document processing and analysis reduces overhead costs. You can automate any repetitive tasks central to your operations and overcome bottlenecks, eliminating costs that arise from manual data entry and processing. You can leverage IDP to boost productivity and streamline workflows across your business operations.

- **Customer satisfaction:** With IDP, you can handle customer documents faster. You can use IDP to automate tasks such as customer onboarding, bookings, and payments that involve documentation. Chatbots can use data from customer documents to respond to customer queries in a more personalized manner. Providing answers and services to customers more quickly enhances customer relationships.[1]

## 1.2 SUBTASKS

- **TSR (Table Structure Recognition)**
- **NER (Named Entity Recognition)**
- **TS (Text Summarization)**

## 1.3 APPLICATIONS

Following is the list along with a brief description of each use case and are some of the most prominent features and applications of IDP:

1. **Healthcare:** IDP improves the management of healthcare records. The healthcare industry must keep immaculate patient records across every touchpoint with a hospital or medical institution. Healthcare businesses use IDP to extract data from patient records and better organize medical documents. The healthcare insurance industry also uses IDP to verify claims and reduce manual paperwork in this field.
2. **Finance:** The financial sector uses IDP to automate several aspects of expense management and invoice processing. Businesses can streamline expense report generation by extracting data from expenses, forms, and business receipts. Financial departments can manage employee and contractor payments with speed and efficiency. For example, an IDP solution can extract figures from financial documents and process data for future payments.
  - I. **Legal:** Businesses in the legal sector can use IDP to analyze contracts. Legal teams use natural language processing (NLP) to analyze a legal contract's terms and obligations. They can extract data from legal documents and court records to build more robust legal cases.
  - II. **Logistics:** Businesses that work in logistics need to track shipments, transit permits, and other vital documents. Companies use IDP for processing documents to reduce the chance of a human error creating a critical mistake. IDP helps with data extraction, validation, and classification, so companies in the logistics sector can speed up logistic functions.
  - III. **Human resources:** Human resources (HR) agents use IDP to extract important information from a candidate's resume. An IDP system saves time and ensures that HR teams focus on choosing between top candidates. The HR industry also uses IDP when managing payroll, leave allotment and other HR functions.

## CHAPTER 2: RESEARCH UNDERWENT

### 2.1 LITERATURE REVIEWS

#### Literature Review of Named Entity Recognition

##### Introduction

Many natural language tasks involve entities, e.g., relation classification, entity typing, named entity recognition (NER), and question answering (QA). Key to solving such entity-related tasks is a model to learn the effective representations of entities. Conventional entity representations assign each entity a fixed embedding vector that stores information regarding the entity in a knowledge base. Although these models capture the rich information in the KB, they require entity linking to represent entities in a text, and cannot represent entities that do not exist in the KB. By contrast, contextualized word representations (CWRs) based on the transformer, such as BERT and RoBERTa provide effective general-purpose word representations trained with unsupervised pretraining tasks based on language modeling. Many recent studies have solved entity-related tasks using the contextualized representations of entities computed based on. However, the architecture of CWRs is not well suited to representing entities for the following two reasons: (1) Because CWRs do not output the span-level representations of entities, they typically need to learn how to compute such representations based on a downstream dataset that is typically small. (2) Many entity-related tasks, e.g., relation classification and QA, involve reasoning about the relationships between entities. Although the transformer can capture the complex relationships between words by relating them to each other multiple times using the self-attention mechanism, it is difficult to perform such reasoning between entities because many entities are split into multiple tokens in the model. Furthermore, the word-based pretraining task of CWRs is not suitable for learning the representations of entities because predicting a masked word given other words in the entity, is clearly easier than predicting the entire entity.

##### LUKE

LUKE is based on a transformer (Vaswani et al., 2017) trained using a large amount of entity-annotated corpus obtained from Wikipedia. An important difference between LUKE and existing CWRs is that it treats not only words, but also entities as independent tokens, and computes intermediate and output representations for all tokens using the transformer (see Figure 1). Since entities are treated as tokens, LUKE can directly model the relationships between entities. LUKE is trained using a new pretraining task, a straightforward extension of BERT's masked language model (MLM) (Devlin et al., 2019). The task involves randomly masking entities by replacing them with [MASK] entities, and trains the model by predicting the originals of these masked entities. We use RoBERTa as base pre-trained model, and conduct pretraining of the model by simultaneously optimizing the objectives of the MLM and our proposed task. When applied to downstream tasks, the resulting model can compute representations of arbitrary entities in the text using [MASK] entities as inputs. Furthermore, if entity annotation is available in the



task, the model can compute entity representations based on the rich entity-centric information encoded in the corresponding entity embeddings. Another key contribution of the paper was that it extends the transformer using entity-aware self-attention mechanism. Unlike existing CWRs, the model needs to deal with two types of tokens, i.e., words and entities. Therefore, it is beneficial to enable the mechanism to easily determine the types of tokens. To this end, we enhance the self-attention mechanism by adopting different query mechanisms based on the attending token and the token attended to.

## **Distilbert-base-cased**

### **Introduction**

The "cased" variant of DistilBERT, known as **distilbert-base-cased**, maintains the original casing of input tokens, preserving distinctions between uppercase and lowercase letters in text. This model achieves efficiency gains through several key optimizations. First, it uses a process called distillation to compress BERT's knowledge into a smaller architecture while maintaining performance. DistilBERT achieves this by removing certain layers and using smaller embeddings during training, resulting in faster training times and lower memory requirements compared to its predecessor.

Despite its reduced size, **distilbert-base-cased** retains strong performance across various NLP tasks, including text classification, question answering, and named entity recognition. It benefits from BERT's pre-training on large-scale corpora like Wikipedia and BookCorpus, which imbue it with a deep understanding of language semantics and context. This pre-training allows the model to generate high-quality embeddings that capture intricate relationships between words and phrases, facilitating accurate and context-aware predictions in downstream applications.

1. **Input Embedding:** Text inputs are tokenized into subword units using the WordPiece tokenizer. These tokens are then converted into numerical embeddings using learned token embeddings and positional embeddings.
2. **Transformer Encoder Stacks:** DistilBERT consists of multiple Transformer encoder layers. Each layer includes a self-attention mechanism and feed-forward neural network. The self-attention mechanism allows the model to weigh the importance of different tokens in the context of each other, capturing dependencies and relationships within the input sequence. This process is computationally efficient due to the reduced number of attention heads and smaller hidden layer sizes compared to BERT.

3. **Layer Distillation:** During training, DistilBERT employs a process known as knowledge distillation. It learns to mimic the behavior of a larger, more complex model (like BERT) by minimizing the difference between its predictions and those of the larger model. This distillation process helps compress BERT's knowledge into a smaller model architecture, reducing memory and computational requirements while maintaining performance.
4. **Output Prediction:** The final layer of DistilBERT outputs logits for various tasks, such as classification or sequence labeling. These logits are then converted into probabilities through a softmax function, producing predictions for specific tasks based on the input text.
5. **Training and Fine-Tuning:** DistilBERT is pretrained on large-scale corpora using unsupervised learning objectives like masked language modeling and next sentence prediction. It can be further fine-tuned on domain-specific labeled data to improve performance on specific downstream tasks, adapting its embeddings and parameters accordingly.

## Literature Review of Table Structure Recognition (TSR)

### Overview

Table Structure Recognition (TSR) is a critical task in document processing, aiming to accurately identify and extract tabular data from various document formats. This involves detecting the presence of tables, recognizing their structure, and extracting the content of individual cells. Recent advancements in machine learning and deep learning have significantly improved the accuracy and efficiency of TSR systems. In this review, we discuss key methodologies and algorithms employed in TSR, particularly focusing on the YOLOv8 architecture for table detection, the SLANet algorithm for table structure recognition, and insights from the PubTables1M dataset.

### YOLOv8 Architecture for Table Detection

**Overview:** YOLO (You Only Look Once) is a state-of-the-art object detection framework known for its speed and accuracy. YOLOv8 represents one of the latest iterations in this series, incorporating several enhancements over its predecessors. It is designed to perform object detection in real-time, making it suitable for applications requiring quick and accurate detection.

### Key Features:

- **Single-Stage Detector:** Unlike traditional object detectors that use a two-stage process (region proposal and classification), YOLOv8 is a single-stage detector. This approach significantly speeds up the detection process by predicting bounding boxes and class probabilities directly from the input image in a single pass.
- **Backbone Network Enhancements:** YOLOv8 employs an improved backbone network with deeper and more efficient layers, enhancing its feature extraction capabilities. This allows the model to capture more complex patterns and finer details in images, which is crucial for accurately detecting tables with varied layouts and structures.
- **Anchor Boxes and Grid Cells:** The model uses anchor boxes and grid cells to predict object locations and dimensions. In the context of table detection, this helps in accurately localizing table boundaries within a document image.
- **Feature Pyramid Network (FPN):** YOLOv8 integrates an FPN to improve multi-scale detection. This is particularly useful for table detection as tables can vary significantly in size within a document.

Application in Table Detection: The fine-tuning of YOLOv8 for table detection involves training the model on a dataset specifically curated for tables. This dataset includes diverse examples of tables from various document types and formats, allowing the model to generalize well. The high accuracy of YOLOv8, combined with its speed, makes it an ideal choice for the table detection component of the TSR pipeline.

### SLANet Algorithm for Table Structure Recognition

**Overview:** SLANet (Structure Line Adjustment Network) is an advanced deep learning algorithm designed specifically for table structure recognition. It addresses the challenges of accurately identifying the hierarchical and functional relationships within tables, which is essential for extracting meaningful data.

### Key Features:

- **Line Adjustment Mechanism:** SLANet employs a unique line adjustment mechanism to accurately identify the boundaries of rows, columns, and merged cells within a table. This mechanism helps in aligning detected text regions with the actual table structure, ensuring precise boundary detection.

- **Hierarchical Structure Recognition:** The algorithm is capable of recognizing the hierarchical structure of tables, including headers, sub-headers, and data cells. This involves understanding the spatial relationships between different table elements and preserving their logical associations.
- **End-to-End Learning:** SLANet is trained end-to-end, meaning it simultaneously learns to detect text regions and recognize table structure. This integrated approach improves the overall accuracy and efficiency of the table structure recognition process.
- **Robustness to Variations:** SLANet is designed to handle various table layouts, styles, and document conditions. This includes dealing with skewed tables, varying font sizes, and different formatting styles.

**Application in Table Structure Recognition:** SLANet is applied after text detection and recognition have been performed on the detected table images. The recognized text regions are input into SLANet, which then adjusts the boundaries and identifies the structure of the table. The algorithm ensures that the extracted data accurately reflects the original table's layout and logical relationships, making it suitable for further processing and analysis.

## **PubTables1M Dataset**

**Overview:** The PubTables1M dataset is a significant contribution to the field of Table Structure Recognition (TSR), providing the largest dataset specifically designed for TSR research. This dataset includes over one million tables extracted from PubMed Central (PMC) open access articles, offering a diverse and comprehensive resource for training and evaluating TSR models.

### **Key Features:**

- **Diverse Table Styles:** The dataset encompasses a wide variety of table styles, layouts, and structures, reflecting the diversity found in scientific literature. This includes simple tables, complex tables with merged cells, multi-level headers, and tables with various formatting styles.
- **High-Quality Annotations:** Each table in the dataset is annotated with detailed structural information, including the locations of rows, columns, and cells, as well as the hierarchical relationships between table elements. This provides a rich source of training data for developing robust TSR models.
- **Large Scale:** With over one million tables, the PubTables1M dataset offers a substantial volume of data for training deep learning models, enabling them to learn complex patterns and generalize effectively to new table formats.

## **Literature Review on Text Summarization**

### **Introduction:**

In an era where the volume of information is growing exponentially, the need for efficient and effective methods to digest large amounts of text has become increasingly important. Text summarization, the process of condensing a long document into a shorter version while retaining its essential information, has emerged as a crucial technology to address this challenge. This technique not only saves time for readers but also enhances the accessibility and comprehensibility of textual content.

### **Overview of Text Summarization**

Text summarization distills the core ideas and key information from a lengthy document into a concise and coherent summary. This process involves several intricate steps, each contributing to the generation of a high-quality summary. The primary goal is to ensure that the summary accurately reflects the original content while being significantly shorter in length. Text summarization can be broadly categorized into two types: extractive and abstractive.

#### **Extractive Summarization**

Extractive summarization involves selecting significant sentences, phrases, or sections directly from the original text and concatenating them to form a summary. This method relies on identifying the most relevant parts of the text and reordering them if necessary. While extractive summarization maintains the original wording and phrasing, it may sometimes result in summaries that lack coherence and fluency due to the disjointed nature of the selected excerpts.

#### **Abstractive Summarization**

In contrast, abstractive summarization generates summaries by interpreting and paraphrasing the main ideas of the source text. This approach involves understanding the content at a deeper level and producing new sentences that convey the same meaning as the original text. Abstractive summarization is more challenging than extractive summarization as it requires advanced natural language processing (NLP) techniques to ensure that the generated summaries are both accurate and readable.

## Steps in Text Summarization

The process of text summarization can be broken down into several key steps, each playing a vital role in the creation of an effective summary. These steps include:

1. **Input Text:** The initial step in text summarization is the acquisition of the input text. This text can come from various sources such as articles, reports, books, or any other form of written content. The quality and structure of the input text can significantly influence the outcome of the summarization process.
2. **Tokenization:** Tokenization is the process of breaking down the input text into smaller units such as words, phrases, or sentences. This step is crucial as it transforms the text into a format that can be easily processed by NLP algorithms. Tokenization helps in identifying the individual components of the text and facilitates subsequent analysis and processing.
3. **Pre-training:** Pre-training involves training a language model on a large corpus of text to learn the underlying patterns and structures of the language. This step equips the model with a broad understanding of grammar, syntax, and semantics, which is essential for generating coherent and contextually accurate summaries. Pre-trained models such as GPT (Generative Pre-trained Transformer) and BERT (Bidirectional Encoder Representations from Transformers) have proven to be highly effective in various NLP tasks, including text summarization.
4. **Fine-tuning:** Fine-tuning is a critical step that involves adapting the pre-trained model to the specific task of text summarization. During this phase, the model is trained on a smaller, task-specific dataset to refine its understanding and improve its performance. Fine-tuning helps the model to learn the nuances of summarization, such as identifying key information, maintaining coherence, and generating concise summaries.
5. **Summary Generation:** Once the model has been fine-tuned, it can be used to generate summaries from the input text. This step involves applying the trained model to the text, which produces a condensed version that captures the main ideas and key information. The quality of the generated summary depends on the effectiveness of the previous steps and the capability of the model to understand and interpret the content accurately.
6. **Output:** The final step in the text summarization process is the presentation of the generated summary. The output should be a coherent, concise, and informative version of the original text, providing readers with a clear understanding of the main points without the need to read the entire document. The output summary can be used in various applications, including news aggregation, academic research, content curation, and more.

**Algorithms: TextRank**, introduced by Rada Mihalcea and Paul Tarau in 2004, is inspired by the PageRank algorithm used by Google to rank web pages. TextRank builds a graph where sentences or words are nodes, and edges represent the similarity or relationship between nodes.

The algorithm iteratively calculates the importance of each node based on its connections, ultimately identifying the most significant sentences or keywords for summarization.

**TF-IDF (Term Frequency-Inverse Document Frequency)** is a widely used technique for information retrieval and text summarization. It evaluates the importance of a term within a document relative to a corpus, which aids in identifying key terms for summarization purposes. This section outlines the working principles of TF-IDF, its advantages and disadvantages in text summarization, the results of applying TF-IDF in a summarization task, and the deployment of this technique on AWS SageMaker.

**BART**, introduced by Facebook AI, is a versatile model that combines the strengths of bidirectional and auto-regressive transformers. BART is pre-trained using a denoising autoencoder approach, where the model is trained to reconstruct original text from corrupted input. This pre-training involves various types of text corruption such as token masking, token deletion, and sentence permutation, enabling BART to understand and generate natural language effectively. BART excels in abstractive summarization due to its ability to generate high-quality summaries that are not merely extracts of the input text but rephrased and synthesized versions that capture the essence of the original content. Its architecture allows it to model long-range dependencies and context, making it particularly powerful for generating coherent and fluent summaries.

**PEGASUS**, developed by Google Research, introduces a novel pre-training strategy specifically tailored for abstractive text summarization. The key innovation of PEGASUS lies in its pre-training objective, where the model learns to generate missing sentences from the input document. These "gap sentences" are strategically extracted based on their importance and informativeness, simulating the task of summarization during pre-training. PEGASUS is trained to predict these missing sentences given the rest of the document, which encourages the model to understand and generate summaries that encapsulate the main points of the text. This pre-training method allows PEGASUS to achieve state-of-the-art performance on various summarization benchmarks, demonstrating its ability to produce concise and informative summaries that go beyond simple extraction and involve sophisticated abstraction and rephrasing.

## CHAPTER 3: PROJECT DEVELOPMENT

### 3.1 NER TASK

Named Entity Recognition (NER) is a key process in **Natural Language Processing (NLP)** that identifies and classifies key elements (entities) in text into predefined categories such as names of persons, organizations, locations, times, quantities, monetary values, and percentages. In the context of **Intelligent Document Processing (IDP)**, NER is particularly valuable for extracting relevant information from unstructured documents. By applying NER, IDP systems can classify documents based on their content, identify and extract specific data points like names, dates, and amounts, and organize this extracted data for further processing or analysis. This capability enables the automation of workflows that involve large volumes of documents, such as invoice processing, legal document review, and customer service tasks. Integrating NER into IDP systems enhances efficiency, reduces manual effort, and improves accuracy in data handling.

#### **Problems:**

While Named Entity Recognition (NER) is a powerful tool in Natural Language Processing (NLP) and Intelligent Document Processing (IDP), it faces several challenges. One major problem is the ambiguity of language; words can have multiple meanings depending on context, making accurate entity classification difficult. Additionally, the diversity of language, including variations in names, formats, and terminologies across different domains, can hinder the effectiveness of NER systems. NER models also struggle with handling rare or previously unseen entities, which can result in incomplete or incorrect data extraction. Moreover, integrating NER into existing systems requires significant computational resources and expertise, and maintaining these systems involves continuous updates to accommodate new data and evolving language patterns. These challenges highlight the complexity of achieving high accuracy and reliability in NER applications within IDP.

**Types:** Named Entity Recognition (NER) can be categorized into various types based on different criteria, enhancing its applicability across diverse contexts. One classification is based on entity categories, distinguishing between standard NER, which identifies entities like persons, organizations, and locations, and fine-grained NER, which delves into more specific subcategories within these standard types.



NER can also be differentiated by the level of supervision used in training: supervised NER relies on annotated data, semi-supervised NER combines limited annotated data with a larger unannotated dataset, and unsupervised NER does not use annotated data, instead identifying entities through patterns in the text. Additionally, NER approaches can be rule-based, using predefined linguistic rules; statistical, employing machine learning algorithms; or neural, leveraging deep learning models for more complex pattern recognition. NER systems can also operate on different levels of text granularity, such as token-based, span-based, or character-based recognition. Furthermore, NER can be designed as general-purpose, applicable to a wide range of texts, or domain-specific, tailored to particular fields like medical, legal, or financial documents.

### **Steps in NER:**

The process of Named Entity Recognition (NER) is essential for extracting structured information from unstructured text, making it a cornerstone of many natural language processing (NLP) applications. The steps involved in NER are comprehensive and systematic, ensuring the accurate identification and classification of entities such as names, locations, dates, and more. Here, we delve into a detailed exploration of each step in the NER process, complemented by highlighting the critical aspects of each stage.

### **Input Text**

The initial step in ner is the acquisition of the input text. This text can come from various sources such as articles, reports, books, or any other form of written content. The quality and structure of the input text can significantly influence the outcome of classifying the named entities.

**Text Preprocessing:** Named Entity Recognition (NER) begins with text preprocessing, a critical step to prepare raw text for analysis. Tokenization is the initial process, breaking down text into individual tokens like words, punctuation marks, and symbols. This step ensures that the text is segmented into manageable units for further processing. Following tokenization, normalization standardizes the text by converting all characters to lowercase, removing punctuation, and handling other textual variations. These actions reduce the complexity of the text and facilitate consistent analysis across different inputs. Proper preprocessing lays the foundation for accurate entity recognition by ensuring uniformity in how the text is processed and interpreted by subsequent NER components.

**Part-of-Speech Tagging:** Part-of-speech tagging is a crucial component of Named Entity Recognition (NER) that assigns grammatical labels to each token in a text corpus. This process categorizes tokens into parts of speech such as nouns, verbs, adjectives, and adverbs, providing linguistic context that aids in identifying entities. By understanding the grammatical roles of tokens within sentences, NER systems can distinguish between different types of entities based on their syntactic behavior. For instance, proper nouns are more likely to be entities like names of people or organizations, while common nouns may not signify entities of interest. Part-of-speech tagging thus enhances the precision of NER by incorporating syntactic information that guides subsequent entity recognition steps. Modern NER systems often integrate advanced tagging models trained on large annotated datasets, leveraging machine learning techniques to achieve high accuracy in part-of-speech assignment. This step is foundational for extracting meaningful entities from text and plays a pivotal role in the overall accuracy and performance of NER systems.

**Feature Extraction:** Feature extraction in Named Entity Recognition (NER) involves identifying and capturing relevant characteristics of tokens within text data to facilitate accurate entity identification and classification. These features encompass lexical, syntactic, and semantic aspects that provide valuable context for understanding the meaning and role of tokens in sentences. Lexical features include attributes such as word morphology (prefixes, suffixes), capitalization, and part-of-speech tags, which help distinguish between entities and non-entities based on their linguistic properties. Syntactic features utilize information about sentence structure and grammatical relationships between tokens to infer entity boundaries and relationships within text. Semantic features leverage word embeddings or contextual embeddings derived from pre-trained language models like BERT, capturing nuanced meanings and contextual associations of tokens beyond their surface forms. By integrating these diverse features, NER systems enhance their ability to recognize and classify entities accurately across different domains and languages. Feature extraction is a fundamental step that bridges raw text data with advanced machine learning models, enabling robust entity recognition capabilities essential for various applications in natural language processing, such as information extraction, document categorization, and sentiment analysis.

**Model Training:** In Named Entity Recognition (NER), model training involves leveraging annotated datasets where entities in text are pre-labeled for machine learning or deep learning models. Supervised learning approaches, such as Conditional Random Fields (CRFs), Long Short-Term Memory (LSTM) networks, or Transformer-based models like BERT, are commonly used. These models learn to recognize patterns and features associated with different types of entities from the labeled data. During training, the model adjusts its parameters based on the examples provided, optimizing its ability to generalize and accurately identify entities in new, unseen text. The training process involves iterative adjustments and evaluations to improve performance metrics such as precision, recall, and F1 score. Model selection depends on the specific requirements of the NER task, considering factors like dataset size, domain specificity, and computational resources. Once trained, the model becomes capable of autonomously identifying entities based on learned patterns and features, forming a crucial component of automated information extraction systems. Effective model training is essential for achieving high accuracy and reliability in NER applications across diverse domains.

**Entity Recognition:** Entity recognition is the core operational phase in Named Entity Recognition (NER) where the trained model is applied to new, unseen text to identify and classify entities based on learned patterns and features. This step involves processing input text through the model, which outputs predictions of entity labels (e.g., person, organization, location) and their corresponding spans within the text. The model utilizes its knowledge acquired during training to analyze each token or sequence of tokens in the text and assign appropriate entity labels based on the probability distribution learned from the training data. Entity recognition is typically performed using efficient algorithms that optimize computational resources while ensuring accurate and real-time processing of text data. The output of this step includes structured data representations of recognized entities, facilitating downstream tasks such as information retrieval, database population, and content analysis. The accuracy and performance of entity recognition directly impact the overall effectiveness of NER systems, influencing their ability to extract meaningful insights from unstructured textual information across various applications and domains. Effective entity recognition capabilities are fundamental to enhancing automation, scalability, and accuracy in natural language processing tasks, supporting advanced functionalities in areas such as healthcare, finance, and customer service.

## **Introduction:**

**LUKE:** Named Entity Recognition (NER) is a pivotal task in natural language processing (NLP) aimed at identifying and categorizing named entities in text.

Leveraging state-of-the-art models like LUKE (Language-Understanding Knowledge-Equipped models), which integrates knowledge graphs into its architecture, enhances NER by incorporating external knowledge for entity disambiguation and contextual understanding. LUKE, based on the Transformer architecture, offers robust capabilities in capturing complex relationships and semantic meanings from text data, crucial for accurate entity recognition across various domains.

**CoNLL-2003**, a benchmark dataset for NER, provides annotated examples in English across multiple entity types (e.g., persons, organizations, locations) essential for training and evaluating NER models. Integrating LUKE with Hugging Face Transformers, a popular library for NLP, facilitates seamless implementation and fine-tuning of pre-trained models on CoNLL-2003 data. This combination harnesses Transformer-based architectures' power to handle long-range dependencies and contextual information effectively, enhancing NER performance compared to traditional methods.

### **LUKE: Algorithm**

LUKE (Language-Understanding Knowledge-Equipped) model integrates knowledge graphs into Transformer-based architectures for advanced Named Entity Recognition (NER) tasks, such as those using the CoNLL-2003 dataset and implemented through Hugging Face Transformers. The algorithm begins with tokenization, breaking text into manageable units, followed by encoding these tokens using Transformer layers to capture contextual relationships and semantic meanings. LUKE enhances this process by incorporating structured knowledge from external sources, facilitating entity disambiguation and enriching contextual understanding beyond textual patterns alone.

During training, LUKE fine-tunes its parameters on labeled datasets like CoNLL-2003, optimizing its ability to recognize and classify entities such as persons, organizations, and locations. This training involves minimizing a loss function that compares predicted entity labels with ground truth annotations, ensuring the model learns to generalize from examples effectively. Post-training, the model undergoes evaluation on validation sets to assess its performance metrics, including precision, recall, and F1 score, crucial for benchmarking against state-of-the-art NER systems.

## Steps in LUKE

1. **Environment Setup and Dependency Management:** Initially, the script ensures all necessary Python libraries and dependencies are installed and configured. This includes installing `sequeval` for evaluating sequence labeling tasks and transformers for seamless integration of the LUKE model into the workflow. These libraries are essential for handling evaluation metrics and facilitating interactions with the pre-trained LUKE model.
2. **Model Loading and Configuration:** The LUKE model checkpoint is obtained from the Hugging Face model hub. This model is specifically trained on the CoNLL-2003 dataset, which provides annotated examples of named entities in English text. Once loaded, the model is configured to leverage GPU acceleration, optimizing computational performance and allowing efficient processing of large-scale data during inference.
3. **Dataset Preparation:** The dataset (`eng.testb`) undergoes preprocessing to organize the data into structured documents. Each document contains words, corresponding entity labels (such as persons, organizations, and locations), and delineates sentence boundaries. This preparation step ensures that the dataset is formatted appropriately for subsequent tokenization and model input.
4. **Tokenization and Input Encoding:** The LUKE tokenizer is employed to tokenize each document, converting raw text into numerical representations suitable for model input. This process includes handling subword segmentation to capture nuanced linguistic features, which is critical for accurate and context-aware entity recognition. By encoding text inputs into tokenized sequences, the model can effectively analyze and classify entities based on their surrounding context and semantic meaning.
5. **Model Inference and Evaluation:** For each batch of processed examples, LUKE applies its span-based entity classification capabilities (`LukeForEntitySpanClassification`). This model architecture allows LUKE to predict entity labels for token spans within the input text, leveraging its knowledge-enhanced embeddings and contextual understanding derived from extensive pre-training on large-scale corpora. Evaluation metrics, such as precision, recall, and F1 score, are computed using `sequeval.metrics.classification_report` to assess the model's performance and accuracy in identifying various types of entities.

6. **Result Analysis and Interpretation:** The model's predictions are analyzed in detail and compared against ground truth labels from the dataset. This analysis provides insights into LUKE's effectiveness in NER tasks, highlighting its ability to discern and classify entities accurately across different domains and linguistic contexts. By integrating LUKE with Hugging Face Transformers, the workflow exemplifies advanced capabilities in natural language processing, demonstrating how state-of-the-art models can be leveraged to enhance entity recognition tasks through sophisticated language understanding and contextual modeling.

### **Advantages of LUKE**

1. **Integration of Knowledge Graphs:** LUKE incorporates structured knowledge graphs into its architecture, enhancing its ability to contextualize and disambiguate entities based on external knowledge. This integration enriches entity recognition by leveraging comprehensive information about relationships and attributes associated with entities.
2. **Transformer Architecture:** Built upon the Transformer architecture, LUKE benefits from its capability to capture long-range dependencies and contextual information effectively. This architecture enables LUKE to understand text at a deeper semantic level, improving the accuracy and robustness of entity recognition tasks.
3. **Contextual Embeddings:** LUKE utilizes knowledge-enhanced embeddings derived from extensive pre-training on large-scale datasets. These embeddings encode rich semantic information and linguistic patterns, facilitating superior performance in NER by understanding nuances in context and entity relationships.
4. **Generalization Across Domains:** Due to its pre-training on diverse datasets, LUKE demonstrates strong generalization across different domains and languages. This capability allows it to adapt well to varying linguistic styles and domain-specific terminologies, making it versatile for applications beyond NER, such as question answering and document understanding.
5. **Enhanced Accuracy and Efficiency:** LUKE's integration of knowledge graphs and Transformer-based architecture enhances both accuracy and computational efficiency in NER tasks. By effectively managing complex information and optimizing inference processes, LUKE provides reliable and scalable performance in real-world applications.

Overall, LUKE represents a significant advancement in NLP by combining knowledge graphs with state-of-the-art Transformer models, offering enhanced contextual understanding, adaptability across domains, and robust performance in entity recognition and related tasks.

## Disadvantages of LUKE

1. **Computational Resources:** LUKE's utilization of large-scale knowledge graphs and Transformer architecture demands significant computational resources, including high memory and processing power. This can limit its deployment on resource-constrained devices or environments without adequate infrastructure.
2. **Fine-Tuning Complexity:** Fine-tuning LUKE for specific tasks like Named Entity Recognition (NER) requires substantial expertise and computational effort. The process involves adjusting model parameters and hyperparameters, which can be time-consuming and challenging for users without extensive NLP expertise.
3. **Knowledge Graph Dependencies:** LUKE's effectiveness heavily relies on the quality and coverage of its underlying knowledge graphs. In scenarios where the knowledge graphs are incomplete or contain inaccuracies, LUKE may struggle to accurately contextualize entities or make incorrect predictions.
4. **Domain Specificity:** Despite its generalization capabilities, LUKE's performance may vary across different domains. Fine-tuning for specific domains or languages is essential to achieve optimal results, and models pretrained on one domain may not seamlessly transfer to another without additional adaptation.
5. **Interpretability and Explainability:** Like many complex deep learning models, LUKE's decision-making process can be opaque and difficult to interpret. Understanding how LUKE arrives at its predictions, especially in complex NLP tasks, remains a challenge, impacting trust and usability in sensitive applications.
6. **Data Dependency:** LUKE's performance heavily depends on the quality, size, and diversity of the training data it is exposed to during pre-training and fine-tuning. Limited or biased training data can lead to suboptimal performance or reinforce biases in entity recognition tasks.

Despite these challenges, ongoing research and advancements in model architecture, training techniques, and data quality are continuously addressing these limitations to enhance LUKE's capabilities and applicability in real-world NLP applications.

## Evaluation of Results:

- **LOC (Location):** The model achieves high precision (95.58%) and recall (94.78%), resulting in a balanced F1-score of 95.18%. This suggests that the model accurately identifies and categorizes locations in the text data.

- **MISC (Miscellaneous):** Precision (85.53%) and recall (86.88%) for miscellaneous entities are slightly lower compared to LOC, leading to an F1-score of 86.20%. There is a slight trade-off between precision and recall, indicating room for improvement in correctly classifying miscellaneous entities.
- **ORG (Organization):** The model shows high precision (92.87%) and recall (94.96%) for organizations, resulting in a robust F1-score of 93.91%. This indicates strong performance in identifying organizational entities.
- **PER (Person):** Precision (96.83%) and recall (97.19%) for persons are notably high, resulting in an impressive F1-score of 97.01%. The model excels in accurately recognizing and categorizing individual names.

### Overall Performance:

- **Micro Average:** The micro average F1-score (94.20%) considers the total number of true positives, false positives, and false negatives across all entity types. It provides an overall measure of the model's accuracy in NER across the entire dataset.
- **Macro Average:** The macro average F1-score (93.07%) calculates the average F1-score across all entity types, treating each class equally. It provides insights into the model's performance without considering class imbalance.
- **Weighted Average:** The weighted average F1-score (94.21%) considers the support (number of true instances) for each class, providing a weighted average that reflects the overall performance across all entity types.

These evaluation metrics collectively indicate that the NER model performs well, particularly in identifying persons and locations, while also showing strong performance for organizations.

### Introduction: DistilBERT

DistilBERT, derived from BERT (Bidirectional Encoder Representations from Transformers), is a compact yet powerful transformer-based model designed by researchers at Hugging Face and Google. It retains much of BERT's architecture while significantly reducing the number of parameters and computational resources required for training and inference.



The "cased" variant of DistilBERT, known as **distilbert-base-cased**, maintains the original casing of input tokens, preserving distinctions between uppercase and lowercase letters in text. This model achieves efficiency gains through several key optimizations. First, it uses a process called distillation to compress BERT's knowledge into a smaller architecture while maintaining performance. DistilBERT achieves this by removing certain layers and using smaller embeddings during training, resulting in faster training times and lower memory requirements compared to its predecessor.

Despite its reduced size, **distilbert-base-cased** retains strong performance across various NLP tasks, including text classification, question answering, and named entity recognition. It benefits from BERT's pre-training on large-scale corpora like Wikipedia and BookCorpus, which imbue it with a deep understanding of language semantics and context. This pre-training allows the model to generate high-quality embeddings that capture intricate relationships between words and phrases, facilitating accurate and context-aware predictions in downstream applications. It represents a streamlined yet proficient alternative to BERT, offering enhanced efficiency without compromising on the robustness and versatility needed for diverse natural language understanding tasks.

### Algorithm of **distilbert-base-cased**

The algorithm of **distilbert-base-cased** involves several key steps that enable it to perform effectively in natural language processing tasks:

1. **Input Embedding:** Text inputs are tokenized into subword units using the WordPiece tokenizer. These tokens are then converted into numerical embeddings using learned token embeddings and positional embeddings.
2. **Transformer Encoder Stacks:** DistilBERT consists of multiple Transformer encoder layers. Each layer includes a self-attention mechanism and feed-forward neural network. The self-attention mechanism allows the model to weigh the importance of different tokens in the context of each other, capturing dependencies and relationships within the input sequence. This process is computationally efficient due to the reduced number of attention heads and smaller hidden layer sizes compared to BERT.

3. **Layer Distillation:** During training, DistilBERT employs a process known as knowledge distillation. It learns to mimic the behavior of a larger, more complex model (like BERT) by minimizing the difference between its predictions and those of the larger model. This distillation process helps compress BERT's knowledge into a smaller model architecture, reducing memory and computational requirements while maintaining performance.
4. **Output Prediction:** The final layer of DistilBERT outputs logits for various tasks, such as classification or sequence labeling. These logits are then converted into probabilities through a softmax function, producing predictions for specific tasks based on the input text.
5. **Training and Fine-Tuning:** DistilBERT is pretrained on large-scale corpora using unsupervised learning objectives like masked language modeling and next sentence prediction. It can be further fine-tuned on domain-specific labeled data to improve performance on specific downstream tasks, adapting its embeddings and parameters accordingly.

This algorithmic framework allows `distilbert-base-cased` to efficiently process and understand natural language, making it a versatile choice for a wide range of NLP applications with reduced computational overhead compared to its predecessor, BERT.

### Steps in `distilbert-base-cased`

6. **Setting Up Fine-Tuning Environment:** The process begins with preparing the environment for fine-tuning `distilbert-base-cased` on the `wnut2017` dataset. This involves installing necessary libraries and dependencies, such as the `tner` library for handling NER tasks. Logging is configured to monitor training progress and capture important information such as training epochs, batch processing, and model evaluation metrics.
7. **Grid Search Configuration:** A `GridSearcher` instance is initialized to facilitate hyperparameter tuning during the fine-tuning process. Hyperparameters like batch size, learning rate, gradient accumulation steps, and others are specified within the `GridSearcher` configuration. This setup allows for systematic exploration of different parameter combinations to optimize model performance on the NER task.
8. **Dataset Preparation:** The `wnut2017` dataset is specified as the training dataset for fine-tuning. This dataset is crucial as it contains annotated examples of named entities, which serve as the ground truth for training the model. The `GridSearcher` is configured to work with this dataset, ensuring that the model learns to recognize entities such as persons, locations, and organizations as defined in the dataset annotations.

9. **Training Execution:** The fine-tuning process begins, where `distilbert-base-cased` is trained on the `wnut2017` dataset. During training, the model iteratively adjusts its internal parameters to minimize prediction errors and improve accuracy in identifying named entities within text inputs. The specified number of epochs and other training parameters dictate how extensively the model learns from the dataset, balancing between underfitting and overfitting to achieve optimal performance.

10. **Model Evaluation:** After training completes, the best-performing model checkpoint is saved. This model represents the configuration that achieved the highest accuracy or performance metric on a validation dataset, as determined by the `GridSearcher`'s evaluation criteria. The saved model can then be used for inference and evaluation on unseen data, such as a test dataset, to assess its ability to generalize to new examples beyond the training set.

11. **Performance Evaluation:** Finally, the trained model is evaluated using standard evaluation metrics for NER tasks, such as precision, recall, and F1-score. These metrics quantify how well the model identifies named entities compared to the ground truth annotations in the test dataset. The evaluation results provide insights into the model's strengths and areas for improvement, guiding further iterations or adjustments in model architecture and training strategies.

By following these steps systematically, the process ensures that `distilbert-base-cased` is effectively fine-tuned and evaluated for NER tasks on the `wnut2017` dataset, leveraging hyperparameter optimization and rigorous evaluation to achieve robust performance in real-world applications.

## Advantages of `distilbert-base-cased`

1. **Efficiency and Speed:** One of the primary advantages of `distilbert-base-cased` is its computational efficiency. By distilling knowledge from the larger BERT model, DistilBERT significantly reduces the number of parameters while maintaining competitive performance. This reduction translates into faster training times and lower memory requirements, making it more accessible for deployment in resource-constrained environments and applications requiring real-time inference. The streamlined architecture allows DistilBERT to perform complex NLP tasks swiftly without compromising accuracy, thus improving operational efficiency in production settings.

2. **Scalability and Deployment:** DistilBERT's smaller model size enhances scalability across various platforms and devices. It can be easily deployed on edge devices, mobile applications, and cloud-based services due to its reduced computational footprint. This scalability makes DistilBERT suitable for large-scale deployment in diverse domains such as healthcare, finance, and customer service, where real-time processing and responsiveness are critical. Additionally, the model's efficient inference capabilities contribute to cost-effectiveness and scalability in handling large volumes of data and user interactions.
3. **Performance Versatility:** Despite its compact size, `distilbert-base-cased` maintains robust performance across a wide range of natural language processing tasks. It inherits BERT's ability to capture intricate linguistic patterns and semantic relationships, enabling accurate predictions in tasks like text classification, sentiment analysis, and entity recognition. The model's pretrained embeddings effectively generalize across domains, adapting to various languages and text genres with minimal fine-tuning. This versatility allows DistilBERT to excel in both supervised and unsupervised learning scenarios, providing reliable results across different applications and datasets.
4. **Resource Efficiency:** DistilBERT optimizes resource utilization through its efficient utilization of computational resources during training and inference. By employing knowledge distillation, the model achieves a balance between computational efficiency and predictive accuracy, making it a cost-effective solution for organizations and developers. The reduced model size also lowers the carbon footprint associated with training and running NLP models, aligning with sustainability goals in AI research and deployment. Overall, DistilBERT's resource efficiency enhances its practicality and accessibility, enabling widespread adoption in industries requiring scalable and environmentally conscious AI solutions.

### **Disadvantages of `distilbert-base-cased`**

1. **Reduced Performance on Complex Tasks:** While `distilbert-base-cased` achieves efficiency gains by distilling knowledge from BERT, it may sacrifice some performance compared to its larger counterpart. This reduction in model size and parameters can lead to decreased accuracy, especially in tasks requiring nuanced understanding or extensive context modeling. Complex tasks that rely heavily on fine-grained details or long-range dependencies may not be handled as effectively as with larger models like BERT.

2. **Limited Representational Capacity:** Due to its compressed architecture, `distilbert-base-cased` has a smaller representational capacity compared to BERT. This limitation can impact the model's ability to capture and retain intricate semantic nuances and domain-specific knowledge. As a result, it may struggle with tasks that demand in-depth understanding of specialized vocabulary, domain-specific jargon, or uncommon linguistic constructs, potentially leading to suboptimal performance in such scenarios.
3. **Dependency on Pretraining Quality:** The effectiveness of `distilbert-base-cased` heavily relies on the quality and diversity of the pretraining data used to distill knowledge from BERT. Inadequate or biased training data can limit the model's ability to generalize across different domains and linguistic variations, affecting its robustness and adaptability. Moreover, variations in pretraining objectives and data sources can impact the model's performance consistency across different applications and datasets.
4. **Contextual Information Loss:** The distillation process in `distilbert-base-cased` involves simplifying and compressing BERT's knowledge into a smaller model. This compression can lead to a loss of certain contextual information and nuanced representations present in the original BERT model. As a result, the distilled model may struggle to capture subtle contextual cues and dependencies in text, potentially affecting its accuracy in tasks requiring precise understanding of language semantics and relationships.
5. **Limited Support for Specialized Tasks:** While `distilbert-base-cased` performs well across general NLP tasks, it may not offer the same level of support or specialized capabilities as larger, task-specific models or architectures. Tasks requiring domain-specific knowledge, extensive fine-tuning, or specific architectural modifications may not be adequately supported by DistilBERT's generalized framework and reduced parameter count, necessitating alternative solutions for optimal performance.

## Evaluation of Results:

### Corporation

The model achieved a precision of 16%, indicating that when it predicted entities as corporations, only 16% of these predictions were correct. The recall was 14%, suggesting that the model identified 14% of all true corporation entities. The F1-score, which combines precision and recall

into a single metric, was 15% for this category. The support, or the number of true instances of corporations in the test set, was 66.

### **Group**

For the group category, the model showed a precision of 48%, indicating that nearly half of its predictions for group entities were correct. However, the recall was low at 9%, suggesting that the model missed identifying many true group entities. The F1-score for groups was 15%, reflecting a balance between precision and recall. The support for group entities in the test set was 165.

### **Location**

In the location category, the model demonstrated a precision of 51%, indicating that a little over half of its predicted location entities were correct. The recall was 40%, suggesting that the model captured 40% of all true location entities. The F1-score for locations was 45%, indicating a relatively balanced performance between precision and recall. There were 150 true instances of location entities in the test set.

### **Person**

For person entities, the model achieved a high precision of 73%, indicating that it correctly predicted a significant portion of person entities. However, the recall was 44%, suggesting that the model missed identifying some true person entities. The F1-score for persons was 55%, reflecting a strong overall performance in this category. The support for person entities in the test set was 429.

### **Product**

The model's performance on product entities showed a precision of 26%, indicating that only a quarter of its predictions for product entities were correct. The recall was 8%, suggesting that the model identified a small proportion of all true product entities. The F1-score for products was 12%, indicating challenges in both precision and recall. There were 127 true instances of product entities in the test set.

### **Work of Art**

In the work of art category, the model achieved a precision of 31%, indicating that about a third of its predictions for work of art entities were correct. The recall was 11%, suggesting that the model

identified a small proportion of all true work of art entities. The F1-score for works of art was 17%, indicating room for improvement in both precision and recall. The support for work of art entities in the test set was 142.

### **Micro Average**

Overall, across all categories, the model achieved a micro average precision of 54%, recall of 28%, and F1-score of 36%. These metrics reflect the model's aggregate performance across all entity types, weighted by the number of true instances in each category.

### **Macro Average**

The macro average precision was 41%, recall was 21%, and F1-score was 26%. These metrics provide an average performance measure across all categories, without considering class imbalance.

### **Weighted Average**

The weighted average precision was 52%, recall was 28%, and F1-score was 35%. These metrics are weighted by the support of each class, providing an overall performance measure that considers the distribution of true instances across different categories in the test set.

This evaluation report provides insights into how well the NER model performed across various entity types in the test dataset, highlighting areas of strength and areas for improvement in entity recognition tasks.

## 3.2 TSR TASK

Table Structure Recognition (TSR) is a vital component of document processing that focuses on automatically identifying and extracting structured tabular information from unstructured documents. In the realm of Intelligent Document Processing (IDP), TSR plays a crucial role in transforming raw textual data into organized, machine-readable formats.

TSR involves employing advanced algorithms and techniques, often leveraging Optical Character Recognition (OCR), natural language processing (NLP), and machine learning models. These technologies enable systems to detect table boundaries, interpret cell contents, and discern the hierarchical relationships within tables, thereby facilitating efficient data extraction and analysis.

### Project Overview of Table Structure Recognition (TSR)

The Table Structure Recognition (TSR) project aims to develop an advanced system capable of accurately identifying and extracting structured information from tables present in various types of documents. This project is particularly significant for fields that rely heavily on data extraction from reports, invoices, research papers, and other documents where tabular data is prevalent.

#### Core Ideas and Objectives

1. **Automated Detection and Extraction:**

- Develop algorithms to automatically detect tables within documents.
- Extract table boundaries, rows, columns, and individual cells to structure the data

appropriately.

2. **Integration of Advanced Technologies:**

- Utilize Optical Character Recognition (OCR) to convert scanned images of documents into machine-readable text.
- Employ machine learning and deep learning techniques to enhance the accuracy of table detection and structure recognition.
- Implement Natural Language Processing (NLP) for understanding the context and content within tables.

3. **Handling Diverse Document Formats:**

- Ensure the system can process a variety of document formats, including PDFs, images, and text files.



- Address challenges posed by different table layouts, such as complex nested tables, merged cells, and varying table structures.

#### 4. **Scalability and Efficiency:**

- Design the system to handle large volumes of documents efficiently.
- Optimize algorithms for quick processing without compromising accuracy.

#### 5. **User-Friendly Interface:**

- Develop an intuitive user interface that allows users to upload documents, review extracted tables, and correct any recognition errors.
- Provide options for exporting structured data in various formats (e.g., CSV, JSON).

### **Pipeline Generation:**

To achieve this, the TSR task is divided into two main subtasks:

1. **Table Detection**
2. **Functional Analysis and Cell Detection**

#### **Subtask 1: Table Detection**

**Objective:** Develop robust algorithms to accurately detect the presence and boundaries of tables within a document.

#### **Key Activities:**

- **Document Preprocessing:** Convert documents into a standardized format for easier analysis. This might include image enhancement, noise reduction, and text normalization.
- **Table Localization:** Use computer vision techniques to identify regions in the document that contain tables. This involves detecting lines, borders, and whitespace patterns typical of table structures.
- **Algorithm Development:** Implement machine learning models, such as convolutional neural networks (CNNs), to learn and predict the location of tables.
- **Performance Optimization:** Ensure the detection algorithm works efficiently across various document types and formats, handling challenges like skewed tables and varying table styles.

#### **Expected Outcomes:**

- High accuracy in detecting tables in diverse document formats.
- Robustness to different table layouts and document conditions.

#### **Subtask 2: Functional Analysis and Cell Detection**

**Objective:** Analyze the detected tables to accurately identify and extract individual cells and their functional relationships.

**Key Activities:**

- **Cell Boundary Detection:** Develop techniques to identify the boundaries of individual cells within a detected table. This might involve analyzing line intersections, whitespace distribution, and content alignment.
- **Content Classification:** Classify the content of each cell using OCR and NLP techniques to distinguish between headers, data cells, and other functional elements of a table.
- **Hierarchical Structure Recognition:** Understand the hierarchical structure of the table, including row and column relationships, merged cells, and nested tables.
- **Data Extraction:** Extract and structure the content of each cell into a machine-readable format, preserving the logical and functional relationships within the table.

**Expected Outcomes:**

- Accurate detection and classification of table cells.
- Preservation of the functional and hierarchical relationships within the extracted table data.

## Table Detection Methodology

### Objective

The objective of the Table Detection (TD) task is to develop a robust algorithm capable of accurately detecting the presence and boundaries of tables within a variety of document formats. This is achieved by utilizing a YOLOv8 table detector that has been fine-tuned specifically for table data.

### Methodology

#### 1. Document Preprocessing

- **Standardization:** Convert documents into a consistent format suitable for analysis. This might include converting different file types (e.g., PDF, JPEG) into a standard image format.
- **Image Enhancement:** Apply techniques such as noise reduction, contrast adjustment, and binarization to enhance the quality of the document images, making it easier to detect table structures.

#### 2. Table Detection using YOLOv8

- **Model Selection:** Utilize the YOLOv8 (You Only Look Once, version 8) object detection model, known for its real-time detection capabilities and high accuracy.

- **Fine-Tuning:** Fine-tune the YOLOv8 model on a dataset specifically curated for table detection which has been taken in the COCO format and it is taken from the roboflow repository .
- **Data Annotation:** Annotate the training data with bounding boxes around table regions. This helps the model learn the specific features that distinguish tables from other content.
- **Training:** Train the YOLOv8 model using the annotated dataset, adjusting hyperparameters to optimize detection performance. Techniques such as data augmentation (e.g., rotation, scaling, flipping) are used to enhance model robustness.
- **Inference:** Use the fine-tuned YOLOv8 model to predict the locations of tables in new documents. The model outputs bounding boxes that indicate the detected table regions. It generated two labels for both detecting and classifying the tables in bordered or borderless structures.

borderless 0.55

Impacts on participating teachers			
Knowledge/appreciation of school system and education in the partner countries	-0,1505	-0,1636	-0,1349
Foreign language competence	-0,0545	-0,0997	-0,0519
Social skills and personal commitment	-0,2558	-0,2235	-0,1302
Professional knowledge and abilities	-0,2145	-0,2319	-0,1003
Impacts on the school as a whole			
European/International dimension of the school	-0,2438	-0,1945	-0,1030
School climate	-0,2976	-0,1810	-0,1012
Innovation in teaching and school management	-0,2586	-0,2557	-0,0928
Training of teachers	-0,1839	-0,1703	-0,0518
Involvement of external actors in the every day school-life	-0,2346	-0,2343	-0,2237
International mobility of pupils	**	**	-0,0583

\* Significance p = 0,000  
 \*\* No significant correlation

## Functional Analysis and XLSX File Generation

### Objective

The objective of this subtask is to perform functional analysis on detected tables and generate XLSX files containing the structured data. This involves feeding the detected table images to the PaddleOCR toolkit for text detection and recognition, followed by table structure recognition using the SLANet algorithm.

## Methodology

### 1. Input Preparation

- **Table Images:** Use the output from the table detection phase (detected table images) as the input for this phase.
- **Standardization:** Ensure the table images are in a consistent format suitable for processing by PaddleOCR.

### 2. Text Detection and Recognition using PaddleOCR

- **PaddleOCR Toolkit:** Utilize PaddleOCR, an open-source OCR tool, for text detection and recognition within the table images.
  - **Text Detection:** Detect text regions within the table images. PaddleOCR's text detection module identifies the locations of text blocks, providing bounding boxes around each detected text region.
  - **Text Recognition:** Recognize the text within the detected regions. The text recognition module converts the detected text regions into machine-readable text, providing the content of each cell within the table.

### 3. Table Structure Recognition using SLANet

- **SLANet Algorithm:** Implement the SLANet (Structure Line Adjustment Network) algorithm for table structure recognition. SLANet is designed to understand the spatial relationships and hierarchical structure within tables.
  - **Structure Line Adjustment:** Adjust the detected text regions to align with the actual table structure. This involves identifying and correcting the boundaries of rows, columns, and merged cells.
  - **Hierarchical Relationships:** Recognize the hierarchical relationships within the table, including headers, sub-headers, and data cells. SLANet processes the spatial arrangement of text blocks to accurately represent the table's structure.

### 4. Data Extraction and Structuring

- **Cell Content Extraction:** Extract the content of each cell based on the recognized text and the table structure identified by SLANet.
- **Logical Grouping:** Group the extracted content logically according to the table's hierarchical structure, preserving the relationships between headers and data cells.

### 5. XLSX File Generation

- **Python Libraries:** Utilize Python libraries such as `openpyxl` or `pandas` to generate XLSX files from the structured data.
  - **Workbook Creation:** Create a new Excel workbook and add a worksheet for each detected table.
  - **Data Insertion:** Insert the extracted and structured data into the appropriate cells in the worksheet. Ensure that the layout and formatting of the table are preserved.
  - **Formatting:** Apply necessary formatting to the XLSX file, such as adjusting column widths, setting header styles, and adding borders to cells for better readability.

### 6. Validation and Refinement

- **Manual Review:** Conduct a manual review of the generated XLSX files to ensure the accuracy and completeness of the extracted data.

- **Error Correction:** Identify and correct any errors or inconsistencies in the extracted data or table structure recognition.
- **Automated Validation:** Implement automated validation checks to compare the generated XLSX files against a ground truth dataset, measuring performance metrics such as accuracy and completeness.

B	
力指标可以从宏观上反映一个区域水资源的丰裕/稀缺程度和开发利用程度，却不能反映一个流域的水资源	
评价方法 总体供需平衡比例 管理性缺水 水质性缺水	可选评价指标 人类耗水量占人类可
工程性缺水 资源性缺水 混合性缺水 异常情况下的水资源风险	工程供水能力与水资源量之比 水资
水量保障程度 水质保障程度 水价承受能力	水量的比重 人均农村生活供水量占
水分分配社会公平 水量保障程度 水质保障程度 经济承受能力	企业平均停水时间 有效灌溉面积中
水生态压力 水生态状态 水生态响应的分析和诊断，必须尽可能全面地	三类以下河段比例 累计地下水超采

## Evaluation of the Table Structure Recognition (TSR) Pipeline

The evaluation of the TSR pipeline involves assessing the performance of both the table detection and the functional analysis and cell detection components. The evaluation metrics used for this purpose are the F1 score for table detection and the TEDS (Table Entity Detection Score) for functional analysis and cell detection.

### Table Detection Evaluation

#### Metric: F1 Score

- **F1 Score:** The F1 score is a measure of a test's accuracy, considering both precision and recall. It is the harmonic mean of precision and recall, providing a single metric that balances both false positives and false negatives.

- **Result:** The table detection component of the TSR pipeline achieved an F1 score of 92.66.

#### Analysis:

- **High Precision:** The model correctly identified a high proportion of actual table regions, indicating effective minimization of false positives.
- **High Recall:** The model successfully detected most of the true table regions present in the documents, demonstrating robust identification capabilities.
- **Balanced Performance:** The high F1 score indicates a well-balanced performance in terms of precision and recall, reflecting the reliability of the table detection algorithm.

### Functional Analysis and Cell Detection Evaluation

### **Metric: TEDS (Table Entity Detection Score)**

- **TEDS:** TEDS is a comprehensive metric for evaluating table structure recognition. It considers the accuracy of detecting individual table entities (such as cells, rows, and columns) and the correctness of the overall table structure.
- **Result:** The functional analysis and cell detection component of the TSR pipeline achieved a TEDS of 95.08.

### **Analysis:**

- **Entity Detection:** The high TEDS indicates a high level of accuracy in detecting individual table entities, such as cells and their contents.
- **Structure Recognition:** The score reflects the system's ability to correctly identify the hierarchical relationships and structure within tables, preserving the logical layout and functional associations.
- **Overall Performance:** The high TEDS score signifies the system's robustness in extracting and structuring table data accurately, facilitating reliable data integration and analysis.

### **Combined Evaluation**

The combined evaluation of the TSR pipeline, considering both the F1 score for table detection and the TEDS for functional analysis and cell detection, demonstrates the overall effectiveness and reliability of the system.

## **Deployment and Web Application**

### **Objective**

The objective was to deploy the Table Structure Recognition (TSR) pipeline as a web application, providing users with an accessible and efficient tool to extract table data from PDF documents. This deployment was achieved using Flask and hosted on AWS, ensuring scalability and reliability.

### **Deployment Methodology**

#### **1. Pipeline Deployment using Flask**

- **Flask Framework:** Flask, a lightweight WSGI web application framework in Python, was chosen for its simplicity and flexibility in handling web requests.

- **API Endpoints:** Developed RESTful API endpoints to handle the upload of PDF documents, process the documents through the TSR pipeline, and return the extracted table data.

- **Upload Endpoint:** Allows users to upload PDF files to the server.

- **Processing Endpoint:** Handles the processing of uploaded PDFs, invoking the table detection and functional analysis components of the TSR pipeline.

- **Download Endpoint:** Provides the processed table data, typically in XLSX format, for download.

## 2. AWS Interface

- **AWS Services:** Utilized various AWS services to host and scale the application, ensuring high availability and performance.

- **EC2 Instances:** Deployed the Flask application on AWS EC2 instances, providing the necessary computational resources for running the TSR pipeline.

- **S3 Storage:** Used AWS S3 for storing uploaded PDF documents and generated table images, ensuring secure and scalable storage solutions.

- **Lambda Functions:** Implemented AWS Lambda functions for handling specific tasks within the pipeline, such as triggering processing upon file upload.

- **API Gateway:** Set up AWS API Gateway to manage and route API requests securely.

## 3. Web Application Features

- **User Interface:** Designed a user-friendly web interface allowing users to easily upload PDF documents and view the extracted table data.

- **File Upload:** Simple drag-and-drop interface for uploading PDF files.

- **Processing Status:** Real-time status updates on the processing of uploaded documents.

- **Table Image Display:** Visualization of detected table images, allowing users to verify the detection results.

- **Download Option:** Option to download the extracted table data in XLSX format.

## 4. Processing Flow

- **PDF Upload:** Users upload a PDF document through the web interface.

- **Table Detection:** The uploaded PDF is processed by the table detection component using the fine-tuned YOLOv8 model, identifying and extracting table regions.

- **Text Recognition:** Detected table images are passed to the PaddleOCR toolkit for text detection and recognition.

- **Structure Recognition:** SLANet algorithm is applied to recognize the table structure and extract cell data.

- **XLSX Generation:** The extracted data is structured and formatted into an XLSX file, ready for download by the user.

### **Benefits and Outcomes**

- **Accessibility:** The web application provides an easy-to-use interface for users to extract table data from PDF documents without needing specialized software or technical knowledge.
- **Scalability:** Leveraging AWS infrastructure ensures the application can handle high volumes of requests and large document sizes efficiently.
- **Efficiency:** The deployment of the TSR pipeline as a web service enables quick and accurate extraction of table data, saving time and reducing manual effort.
- **Reliability:** The use of Flask and AWS services ensures a robust and reliable application, capable of operating continuously with minimal downtime.

### **Conclusion**

The deployment of the TSR pipeline using Flask and AWS has resulted in a powerful web application that simplifies the extraction of table data from PDF documents. Users can easily upload PDFs, and the application processes these documents to generate accurate table images and structured data in XLSX format. This deployment showcases the effectiveness of combining advanced machine learning models with scalable cloud infrastructure to provide a valuable tool for data extraction and analysis.

## **3.3 TS TASK**

In an era where the volume of information is growing exponentially, the need for efficient and effective methods to digest large amounts of text has become increasingly important. Text summarization, the process of condensing a long document into a shorter version while retaining its essential information, has emerged as a crucial technology to address this challenge. This technique not only saves time for readers but also enhances the accessibility and comprehensibility of textual content.

### **Overview:**



Text summarization distills the core ideas and key information from a lengthy document into a concise and coherent summary. This process involves several intricate steps, each contributing to the generation of a high-quality summary. The primary goal is to ensure that the summary accurately reflects the original content while being significantly shorter in length. Text summarization can be broadly categorized into two types: extractive and abstractive.

### **Extractive Summarization:**

Extractive summarization involves selecting significant sentences, phrases, or sections directly from the original text and concatenating them to form a summary. This method relies on identifying the most relevant parts of the text and reordering them if necessary. While extractive summarization maintains the original wording and phrasing, it may sometimes result in summaries that lack coherence and fluency due to the disjointed nature of the selected excerpts.

### **Abstractive Summarization:**

In contrast, abstractive summarization generates summaries by interpreting and paraphrasing the main ideas of the source text. This approach involves understanding the content at a deeper level and producing new sentences that convey the same meaning as the original text. Abstractive summarization is more challenging than extractive summarization as it requires advanced natural language processing (NLP) techniques to ensure that the generated summaries are both accurate and readable.

### **Steps in Text Summarization:**

The process of text summarization can be broken down into several key steps, each playing a vital role in the creation of an effective summary. These steps include:

#### **1. Input Text**

The initial step in text summarization is the acquisition of the input text. This text can come from various sources such as articles, reports, books, or any other form of written content. The quality and structure of the input text can significantly influence the outcome of the summarization process.

#### **2. Tokenization**

Tokenization is the process of breaking down the input text into smaller units such as words, phrases, or sentences. This step is crucial as it transforms the text into a format that can be easily processed by NLP algorithms. Tokenization helps in identifying the individual components of the text and facilitates subsequent analysis and processing.

### **3. Pre-training**

Pre-training involves training a language model on a large corpus of text to learn the underlying patterns and structures of the language. This step equips the model with a broad understanding of grammar, syntax, and semantics, which is essential for generating coherent and contextually accurate summaries. Pre-trained models such as GPT (Generative Pre-trained Transformer) and BERT (Bidirectional Encoder Representations from Transformers) have proven to be highly effective in various NLP tasks, including text summarization.

### **4. Fine-tuning**

Fine-tuning is a critical step that involves adapting the pre-trained model to the specific task of text summarization. During this phase, the model is trained on a smaller, task-specific dataset to refine its understanding and improve its performance. Fine-tuning helps the model to learn the nuances of summarization, such as identifying key information, maintaining coherence, and generating concise summaries.

### **5. Summary Generation**

Once the model has been fine-tuned, it can be used to generate summaries from the input text. This step involves applying the trained model to the text, which produces a condensed version that captures the main ideas and key information. The quality of the generated summary depends on the effectiveness of the previous steps and the capability of the model to understand and interpret the content accurately.

### **6. Output**

The final step in the text summarization process is the presentation of the generated summary. The output should be a coherent, concise, and informative version of the original text, providing readers with a clear understanding of the main points without the need to read the entire document. The output summary can be used in various applications, including news aggregation, academic research, content curation, and more.

## Text summarization using TextRank

### Introduction:

Text summarization is an essential task in natural language processing (NLP), enabling the extraction of key information from large volumes of text. One popular algorithm for this purpose is TextRank, a graph-based ranking algorithm that leverages the structure of the text to generate summaries. In this report, we will explore the TextRank algorithm, its advantages and disadvantages, and the results of its application. Additionally, we will discuss the deployment of the TextRank algorithm on AWS.

### TextRank Algorithm:

TextRank, introduced by Rada Mihalcea and Paul Tarau in 2004, is inspired by the PageRank algorithm used by Google to rank web pages. TextRank builds a graph where sentences or words are nodes, and edges represent the similarity or relationship between nodes. The algorithm iteratively calculates the importance of each node based on its connections, ultimately identifying the most significant sentences or keywords for summarization.

### Steps in TextRank for Summarization:

1. **Preprocessing:** Tokenize the text, remove stop words, and perform stemming or lemmatization.
2. **Graph Construction:** Create a graph with sentences as nodes and edges representing sentence similarity based on the cosine similarity of word vectors.
3. **Rank Calculation:** Apply the PageRank algorithm to compute the rank of each sentence.
4. **Sentence Selection:** Select the top-ranked sentences to form the summary.

## Advantages and Disadvantages of TextRank

### Advantages:

**Automatic:** TextRank is an automated process requiring no human intervention, capable of summarizing large texts quickly.

**Unbiased:** The algorithm is unbiased, summarizing text based purely on keyword frequency without considering the author's perspective.

**Time-Saving:** TextRank can provide quick summaries, saving time and effort for users.

**Consistency:** The algorithm produces consistent summaries by following a fixed set of rules.

**Customizable:** TextRank can be customized to prioritize specific keywords or phrases, tailoring the summary to specific needs.

#### **Disadvantages:**

**Limited Context:** TextRank may miss important context not captured by keywords, affecting the summary's comprehensiveness.

**Limited Accuracy:** The algorithm may struggle with poorly written texts or grammatical errors, impacting accuracy.

**Limited Understanding:** Lacking human-like understanding, TextRank may fail to grasp nuances, sarcasm, or irony.

**Limited Coverage:** TextRank is more effective for factual texts and may not perform well with creative or complex texts.

**Limited Creativity:** The algorithm cannot generate creative summaries beyond the scope of the original text.

#### **Results and Evaluation**

The effectiveness of the TextRank algorithm was evaluated using ROUGE and BLEU scores, common metrics for summarization quality assessment.

- **ROUGE Score:**
  - Precision: 1.000
  - Recall: 0.414
  - F1-Score: 0.586
- **BLEU Score:** 0.694

These scores indicate that while TextRank is highly precise, it may miss some relevant information (as indicated by the lower recall), leading to a moderate F1-Score. The BLEU score suggests a relatively high quality of the generated summaries.

## Deployment on AWS

Deploying the TextRank algorithm on AWS involves several steps to ensure scalability and accessibility. Here's a brief overview of the deployment process:

### 1. Setting Up the Environment:

- Launch an EC2 instance with appropriate configurations.
- Install necessary dependencies (Python, NLTK, NumPy, etc.).

### 2. Developing the Application:

- Implement the TextRank algorithm using Python.
- Create an API using Flask to handle requests and return summaries.

### 3. Containerization:

- Use Docker to containerize the application, ensuring consistency across different environments.

### 4. Deploying with AWS Services:

- Push the Docker image to Amazon ECR (Elastic Container Registry).
- Deploy the containerized application using Amazon ECS (Elastic Container Service) or AWS Fargate for serverless deployment.

### 5. Setting Up API Gateway:

- Configure AWS API Gateway to handle incoming requests and route them to the ECS service.

### 6. Monitoring and Maintenance:

- Use AWS CloudWatch to monitor the application's performance and set up alerts for any issues.
- Regularly update the application and dependencies to maintain security and performance.

# TF-IDF for Text Summarization and Deployment on AWS SageMaker

## Introduction

TF-IDF (Term Frequency-Inverse Document Frequency) is a widely used technique for information retrieval and text summarization. It evaluates the importance of a term within a document relative to a corpus, which aids in identifying key terms for summarization purposes. This report outlines the working principles of TF-IDF, its advantages and disadvantages in text summarization, the results of applying TF-IDF in a summarization task, and the deployment of this technique on AWS SageMaker.

## Working Principles of TF-IDF

TF-IDF is calculated as the product of two statistics: Term Frequency (TF) and Inverse Document Frequency (IDF).

**Term Frequency (TF):** This measures how frequently a term appears in a document. It is calculated as the ratio of the number of times the term appears in the document to the total number of terms in the document.

$$TF_{t,d} = \frac{\text{Count of term } t \text{ in document } d}{\text{Total terms in document } d}$$

**Inverse Document Frequency (IDF):** This measures how important a term is. While computing TF, all terms are treated equally important. IDF reduces the weight of terms that appear very frequently across the corpus and increases the weight of terms that appear less frequently.

$$IDF_t = \log \left( \frac{\text{Total number of documents}}{\text{Number of documents containing term } t} \right)$$

**TF-IDF Score:** The TF-IDF score for a term in a document is the product of its TF and IDF values.

$$\text{TF-IDF}_{t,d} = \text{TF}_{t,d} \times \text{IDF}_t$$

## Advantages and Disadvantages of TF-IDF for Text Summarization

### Advantages

1. **Simplicity and Efficiency:** TF-IDF is straightforward to implement and computationally efficient, making it suitable for real-time applications and large datasets.
2. **Term Importance:** It considers both the frequency of terms within a document and across the corpus, helping to identify significant and unique words for summarization.
3. **Customization:** TF-IDF can be tailored to weigh certain terms more heavily based on their relevance to the topic, allowing for more focused and accurate summaries.
4. **Minimal Preprocessing:** It requires minimal preprocessing, making it practical for smaller datasets or simpler NLP tasks.

### Disadvantages

1. **Lack of Context:** TF-IDF evaluates terms independently and does not consider the relationships or context in which they appear, potentially missing out on semantic nuances.
2. **Document Length Sensitivity:** It can be biased towards longer documents that may contain more unique terms, regardless of their actual relevance.
3. **Semantic Meaning:** TF-IDF does not capture the semantic meaning of terms, which can result in summaries missing important concepts.
4. **Equal Term Weighting:** It assumes all terms within a document are equally important, which may not always be accurate in capturing the overall meaning.

### Evaluation Metrics

The performance of the TF-IDF-based summarization was evaluated using ROUGE and BLEU scores:

ROUGE Score:

- Precision: 0.787
- Recall: 0.266
- F1-Score: 0.398
- BLEU Score:0.008

These metrics indicate that while TF-IDF can identify key terms, it might not always produce summaries with high semantic coherence, as reflected in the low BLEU score.

## **Research Papers**

Several research papers have explored the use of TF-IDF for text summarization:

1. "Automatic text summarization using TF-IDF weighting scheme" by R. Wan, D. Zhao, and C. Xu, in Proceedings of the 2010 IEEE International Conference on Intelligent Computing and Intelligent Systems (ICIS).
2. "A comparison study of TF-IDF, LSA and multi-words for text classification" by T. Nasukawa and J. Yi, in Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology (HLT-NAACL).
3. "Extractive summarization using continuous vector space models" by R. Nallapati, B. Zhou, and C. Gulcehre, in Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP).
4. "Text summarization with TF-IDF weighted word embedding" by J. Nam and E. Han, in Proceedings of the 2017 IEEE International Conference on Big Data and Smart Computing (BigComp).

These papers suggest that TF-IDF is effective for extractive summarization, particularly when combined with other techniques like continuous vector space models and word embeddings.

## **Deployment on AWS SageMaker**

Deploying a TF-IDF-based summarization model on AWS SageMaker involves the following steps:



1. **Data Preparation:** Load and preprocess the dataset, ensuring it is suitable for TF-IDF computation.
2. **Model Training:** Implement the TF-IDF algorithm to generate term scores for the documents in the dataset.
3. **Model Evaluation:** Evaluate the summarization performance using metrics like ROUGE and BLEU scores.
4. **Deployment:**
  - I. **Create a SageMaker Notebook Instance:** Use the SageMaker console to set up a notebook instance for developing and testing the TF-IDF model.
  - II. **Train the Model:** Use the notebook to run the TF-IDF algorithm on the dataset and generate summaries.
  - III. **Deploy the Model:** Utilize SageMaker's deployment capabilities to create an endpoint for the TF-IDF model, allowing it to be accessed and used for summarization tasks.
  - IV. **Inference:** Make predictions using the deployed endpoint, generating summaries for new documents.

By leveraging AWS SageMaker, the TF-IDF model can be scaled and integrated into larger applications, providing efficient and accessible summarization services.

## **Text Summarization using BART/BERT and PEGASUS:**

### **Introduction:**

BART (**Bidirectional and Auto-Regressive Transformers**) and PEGASUS (**Pre-training with Extracted Gap-sentences for Abstractive Summarization Sequence-to-sequence models**) are state-of-the-art techniques for text summarization. They leverage deep learning models to generate coherent and contextually rich summaries. This report outlines the working principles of BART and PEGASUS, their advantages and disadvantages in text summarization and the results of applying these models in a summarization task.

### **Working Principles of BART:**

BART combines a bidirectional encoder and an autoregressive decoder, similar to a fusion of BERT and GPT. This hybrid approach allows BART to both understand and generate text effectively.

1. **Pre-training:** BART is pre-trained using a denoising autoencoder approach, where the model learns to reconstruct corrupted input text.

- Corruptions: Text is corrupted using techniques like token masking, token deletion, and sentence permutation.
- Objective: The model is trained to predict the original text from the corrupted input, enhancing its ability to generate coherent text.

2. **Fine-tuning:** After pre-training, BART can be fine-tuned on specific tasks such as summarization, translation, and question answering by training on task-specific datasets.

### **Working Principles of PEGASUS:**

PEGASUS is designed specifically for abstractive summarization. It uses a novel pre-training objective where important sentences are masked and the model is trained to generate these sentences from the rest of the text.

1. **Pre-training:** PEGASUS is pre-trained by masking whole sentences (gap-sentences) that are likely to be good summaries.

- Gap-sentences: Important sentences are identified and masked in the input text.
- Objective: The model learns to generate the masked sentences from the rest of the text, which mimics the process of summarization.

2. **Fine-tuning:** Similar to BART, PEGASUS can be fine-tuned on specific datasets for summarization tasks.

## **Advantages and Disadvantages of BART and PEGASUS for Text Summarization:**

### **Advantages:**

#### **1. Contextual Understanding:**

- BART: Its bidirectional encoder understands the full context of the text.
- PEGASUS: Its pre-training on gap-sentences enables it to generate summaries that capture the main points.

**2. High-Quality Summarization:** Both models produce coherent and fluent summaries that are close to human-generated summaries.

#### **3. Versatility:**

- BART: Can be fine-tuned for multiple NLP tasks beyond summarization.
- PEGASUS: Specially optimized for summarization but can also be used for other tasks with fine-tuning.

**4. Performance:** Both models achieve state-of-the-art performance on multiple summarization benchmarks.

### **Disadvantages:**

**1. Computational Resources:** Both models require significant computational resources for training and inference.

**2. Training Data Dependency:** The quality of the generated summaries depends heavily on the quality and size of the fine-tuning dataset.

**3. Complexity:** The models are complex and require expertise to fine-tune and deploy effectively.

**4. Length Constraints:** There are practical limits to the length of input text that these models can handle effectively.

### **Evaluation Metrics:**

The performance of BART and PEGASUS-based summarization was evaluated using ROUGE and BLEU scores:

**BART - ROUGE Score:**

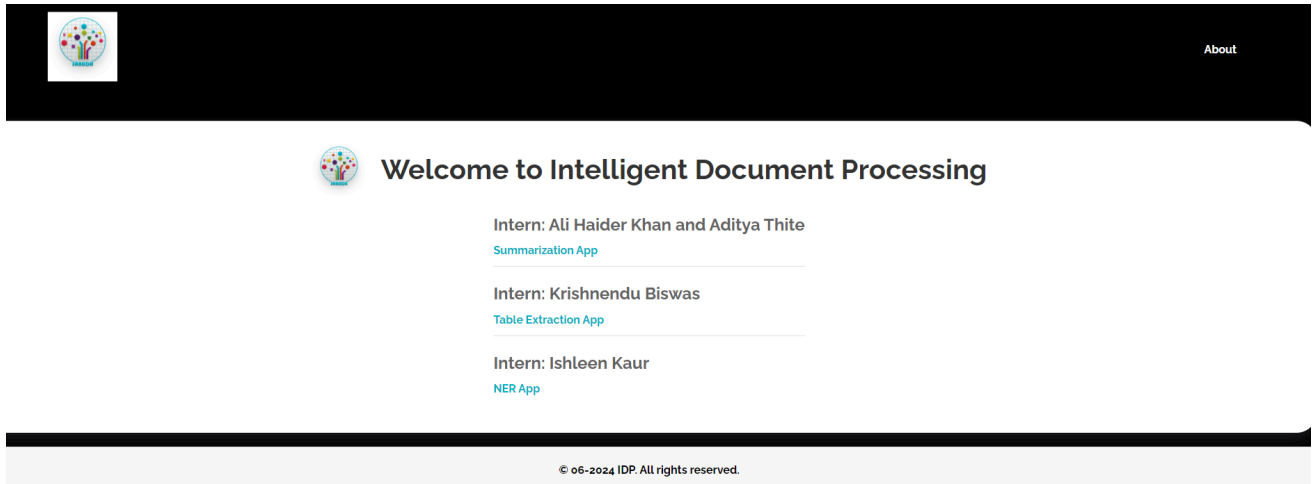
- Precision: 1.00
- Recall: 0.140
- F1-Score: 0.245

**PEGASUS - ROUGE Score:**

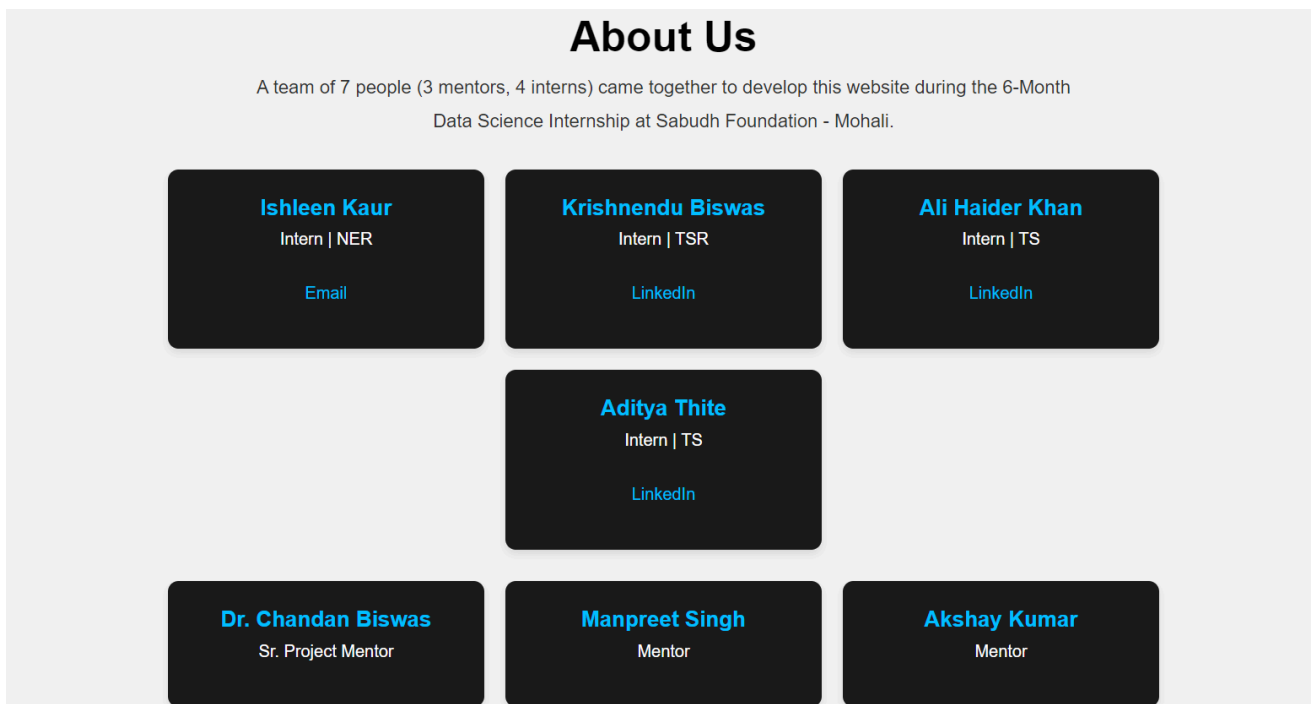
- Precision: 0.969
- Recall: 0.149
- F1-Score: 0.250

## CHAPTER 4: USER INTERFACE

### ‘HOME’ PAGE:



### ‘ABOUT TEAM’ PAGE:



### INDIVIDUAL TASK:

### ‘NER’ PAGE:

The screenshot shows a web interface for Named Entity Recognition. It has a dark background with a central light gray box. At the top of this box, the title "Named Entity Recognition" is displayed in blue. Below the title, the text "Enter text:" is shown. Underneath is a large, empty text input field with a light gray border and a small cursor icon on the right. At the bottom of the central box, there is a black button with the word "Analyze" in white text.

‘TSR’ PAGE:

The screenshot shows a web interface for Table Extractor. It has a dark background with a central light gray box. At the top of this box, the title "Table Extractor" is displayed in white. Below the title, there are two buttons stacked vertically. The top button is green with the text "Choose PDF File" in white. The bottom button is blue with the text "Upload" in white.

‘TS’ PAGE:

## Text Summarization

Enter text to summarize...



BART CNN



Summarize

## CHAPTER 5: CONCLUSION AND FUTURE SCOPE

Intelligent Document Processing (IDP) encompasses technologies aimed at automating document-centric workflows through advanced techniques such as OCR, NLP, and machine learning. Here's an updated version considering IDP:

### Future Scope

#### Named Entity Recognition (NER):

- **Multilingual and Cross-lingual Support:** Enhance NER models within IDP frameworks to process documents in multiple languages effectively, leveraging diverse linguistic datasets and transfer learning techniques.
- **Fine-grained Entity Extraction:** Develop IDP systems capable of identifying detailed entity types across various document types, enhancing semantic understanding and data extraction accuracy.
- **Domain-specific Adaptation:** Tailor NER models integrated into IDP solutions to specific industries like healthcare, legal, and finance, addressing industry-specific terminology and compliance requirements.

#### Text Summarization (TS):

- **Multi-document Summarization:** Advance TS techniques within IDP to summarize information from multiple documents, enabling efficient decision-making and knowledge synthesis across large datasets.
- **Abstractive Summarization:** Develop IDP-driven TS models that generate concise summaries by understanding and synthesizing content, surpassing traditional extractive methods.
- **Personalized Summarization:** Integrate personalized TS capabilities within IDP systems to deliver tailored summaries aligned with user preferences or specific organizational needs, enhancing usability and information delivery.

#### Table Structure Recognition (TSR):



- **Integration with Document Understanding:** Combine TSR capabilities with NLP techniques within IDP frameworks to extract and analyze tabular data in context with surrounding textual content, enabling comprehensive document understanding.
- **Interactive Document Processing:** Develop IDP-driven TSR systems that facilitate interactive exploration and query of tabular data within documents, supporting **dynamic data retrieval and analysis**.
- **Enhanced Accuracy and Robustness:** Continuously improve TSR models integrated into IDP solutions to handle complex table layouts, varying formatting styles, and noisy document conditions, ensuring reliable data extraction and processing.

## Conclusion

Intelligent Document Processing (IDP) solutions leveraging Named Entity Recognition (NER), Text Summarization (TS), and Table Structure Recognition (TSR) are pivotal in transforming document management and data processing:

- NER enhances IDP by accurately identifying and categorizing entities within documents, facilitating precise information retrieval and analysis across diverse document types.
- TS streamlines document consumption within IDP by condensing extensive textual content into concise summaries, preserving critical insights and facilitating informed decision-making.
- TSR automates the extraction and analysis of tabular data within IDP frameworks, improving document understanding and supporting data-driven insights and actions.

As IDP technologies advance, their integration into organizational workflows will drive efficiency, accuracy, and compliance across industries. The future scope includes expanding language support, advancing summarization techniques, and enhancing table recognition capabilities within IDP, positioning NER, TS, and TSR as critical components of intelligent document processing strategies.

## CITATIONS AND REFERENCES

*Website-link:* [https://aws.amazon.com/what-is/intelligent-document-processing/#:~:text=Intelligent%20document%20processing%20\(IDP\)%20is,when%20stock%20levels%20are%20low.d - AWS \(amazon.com\)](https://aws.amazon.com/what-is/intelligent-document-processing/#:~:text=Intelligent%20document%20processing%20(IDP)%20is,when%20stock%20levels%20are%20low.d-AWS%20(amazon.com))

### *NER*

1. <https://arxiv.org/abs/2010.01057>
2. <https://arxiv.org/pdf/2206.12617v3>
3. <https://arxiv.org/pdf/1910.01108>
4. [https://www.researchgate.net/publication/374123033\\_The\\_DistilBERT\\_Model\\_A\\_Promising\\_Approach\\_to\\_Improve\\_Machine\\_Reading\\_Comprehension\\_Models](https://www.researchgate.net/publication/374123033_The_DistilBERT_Model_A_Promising_Approach_to_Improve_Machine_Reading_Comprehension_Models)

### *TSR*

1. <https://universe.roboflow.com/mohamed-traore-2ekkp/table-extraction-pdf>
2. <https://medium.com/@suparnadutta05/table-extraction-drawing-insights-from-table-data-a-survey-report-96c710ebcf55>
3. <https://arxiv.org/abs/2110.00061>