# Survey on Intelligent Document Processing: A Comprehensive Approach to Summarization, NER, Language Conversion, and Plagiarism Detection

Prasanna Avhad[1], Parag Jadhav[1], Sudarshan Madbhavi[1], Ganesh Devnale[1], Prof. Dr. C. A. Ghuge[2]

[1]*Student, Dept. of Artificial Intelligence & Machine Learning Engineering, PES's Modern College of Engineering, Pune, Maharashtra, India*

[2]*Professor, Dept. of Artificial Intelligence & Machine Learning Engineering, PES's Modern College of Engineering, Pune, Maharashtra, India*

---------------------------------------------------------***---------------------------------------------------------

## ABSTRACT

IDP technologies transform the way of information management of various domains with both education and law being on the list. A survey paper encompassing a complete overview of four highly critical modules: PDF/DOCX summarization, named entity recognition, English to Hindi language translation, and plagiarism detection is presented here. IDP technologies enrich academic integrity in the field of education, make easy the processes involved in research, and help bridge language gaps for an inclusive learning environment. From a legal perspective, such technologies would help improve efficiency regarding reviews of documents, contract analysis, and research; thus, professionals would make more certain strides to keep ahead of nonconformities. Some of the major challenges include domain-specific solutions, data quality, ethical concerns, and user-centric design. The paper strongly emphasizes the need for further research and development with integrated, effective, and accessible IDP tools that would meet user requirements over time.

**Keywords**: Intelligent Document Processing, PDF/DOCX Summarization, Named Entity Recognition, Language Conversion, Plagiarism

Detection, Education, Legal Applications, Natural Language Processing, Academic Integrity, Document Review.

## 1. INTRODUCTION

### 1.1 Context & Motivation

Grow of digital information at so frantic a pace that the demand for managing documents has increased much. The volume of information generated by education and the law is so massive, and all that it requires is detailed timely analysis of those documents for the students, educators, lawyers, and researchers. The advent of AI and NLP has just given birth to the Intelligent Document Processing(IDP) concept that could potentially automate the whole task while even improving accuracy and saving human effort. This paper is motivated by an investigation into IDP technologies that summarize, perform named entity recognition (NER), convert languages, and detect plagiarism, all the ways that prove how different processes may be made efficient and scalable in transforming education and legal work landscapes.

### 1.2 Problem Definition

Almost three decades have passed since the journey of document processing began, and significant advancement has been made in designing tools that have overcome the challenges these pose in managing large data, checking for originality, or facilitating

multilingual communication. Starting from the integrity of academics at educational institutions to the long-winded documents one has to labor over to be reviewed by a legal professional, hardly any efficient translation tool is accessible where a multilingual language would be appreciated. A lot of effort still goes into the key details, like picking out people's names and dates. Often, the available solutions fail to tackle all the necessities together.

### 1.3 Goals

The paper reviews IDP technology focusing on the following four major modules.

**PDF/DOCX Summarization** Techniques that automatically summarize documents and help reduce the length of available documents for students to digestible information.

**Named Entity Recognition (NER):** Critical entities like names, organizations, and locations are abstracted from a document to make it easier for information to be retrieved.

**English to Hindi Language Translation Tools:** These are tools that aid frictionless and easy translation of documents, break language barriers, especially in diverse educational and legal backgrounds.

**Plagiarism Detection** Techniques to maintain authenticity and integrity in the academic as well as professional domains.

### 1.4 Modules overview

**Document Summarization** refers to the shortening of long documents into smaller versions without losing main information. Two kinds of methods can be used: either extractive methods involve picking up vital sentences of the text or abstractive methods, where new sentences are produced based on a comprehension of the contents of the document. Recent progress in NLP has greatly improved the quality and coherence of the text summed up through the appearance of transformer-based models, namely BART and GPT.

**Named Entity Recognition (NER):** NER models are designed to identify and classify specific aspects of a document, such as people's names, organization names, places, dates, etc. It is very useful for legal and academic texts because the proper identification of entities helps in proper reading of the content and its context. The latest NER systems are based on the deep learning models, such as BERT and RoBERTa, and on the other transformer-based architectures.

**Language Conversion**, or machine translation, specially in multilingual environments. For example, in countries like India, it becomes very important to have a document that has been translated from one language to the other, yet retaining the exact same meaning and structure as when written originally. Therefore, improvements made to the models in neural machine translation by applying transformers have greatly enhanced translation accuracy as well as fluency, making it feasible for systems to handle complicated documents in multiple languages. For example, the Google Translate API, among others, are widely applied within IDP systems to automatically translate documents in a manner to make them accessible to a wider audience.

**Plagiarism Detection**, especially in the knowledge and publishing industries. Digital content availability on the rise generates a higher risk factor of plagiarism and requires stronger and more efficient systems that are developed specifically for detecting and preventing such cases of plagiarism. Normally, plagiarism detection tools depend on string matching and tokenization and deep learning models that can identify both exact and paraphrased matches. All the systems, upon receiving the submitted documents, compare them with huge database collections of research documents, sites, and so on for similarities to help ensure originality of the content.With all these technologies offered in one IDP platform, a holistic solution to the challenges of digital document processing is provided.

These kinds of systems are most beneficial in industries where volumes of unstructured data are subject to rapid processing with complete accuracy, such as in educational institutions, where students and researchers can benefit from automatic summarization tools that condense lengthy research papers into easy-to-manage summaries. NER can assist lawyers automatically in the identification of significant entities in contracts and case files, thus helping save time and reduce risks due to missing critical information. Language conversion allows for accessible documents from non-native writers; plagiarism detection helps in integrity for academic purpose.

## 2. BACKGROUND

### Theoretical Concepts

Understanding the theoretical foundations of Intelligent Document Processing is fundamental to understand how AI and NLP-driven systems can be more accurately applied in document management. Some of the key concepts involved are:

**2.1 Natural Language Processing (NLP):** NLP is a field of AI whose application is specifically aimed at enabling machines to understand, interpret, and generate human language. It plays an important role in the process of document management, including functionalities such as summarization, translation, and named entity recognition that make the process of document analysis more efficient and scalable.

**2.2 Transformer Models:** High reliance on transformer models in modern NLP is BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative Pre-trained Transformer) models, which are truly strong in the application of attention mechanisms to understand context and apply best to summarization, language translation, and even entity recognition. The ability of transformer models in processing large volumes of text data allows it to be able to sum up, translate texts perfectly, and be able to also extract important entities while solving complex documents.

**2.3 Summarization Techniques:** Summarization can be either extractive or abstractive. Extractive methods mean pulling key sentences directly out of the document, whereas abstractive techniques produce new sentences to convey the meaning of the document. Here, in IDP, a machine summarization technique that automatically condenses long documents, like legal contracts or academic papers, into manageable sections saves time and effort.

**2.4 Named Entity Recognition (NER):** NER is that subset of NLP from where the mainly identified and categorized key entities are extracted from text. Such as names, organizations, dates, and locations. For example, for educational or legal purposes, identifying this in bulk from massive documents can reduce the load on workload and improve the accuracy of information retrieval.

**2.5 Machine Translation:** This automatically translates text from one language to another, using Artificial Intelligence with the help of Text Machine Translation. In IDP, the tools which support translating documents from English into the Hindi Language or any other can provide wider access, especially in multilingual environments like education and legal sectors.

**2.6 Plagiarism Detection:** The plagiarism-detection tools use a combination of algorithms for text matching and apply machine learning in order to identify the unoriginal content presented in a document. Such systems are essential both for integrity in an academic setting and for law enforcement purposes to ensure that theft or abuse of intellectual property does not occur.

## 3. PDF/DOCX SUMMARIZATION TECHNIQUES

Automatic summarization is a vital portion of the IDP systems for the simple reason of an all-inclusive reduction in length without compromising the content. All fields have seen a rapid growth in textual data, making it essential to sum up lengthy documents professionally with efficient summarization techniques, especially in fields such as law, education, and business. Summarization can be broadly categorized into two approaches or types of summarization: extractive and abstractive summarization.

**3.1 Extractive Summarization** is pulling out the most relevant sentences or phrases exactly from the source document. The summary is obtained based on ranking of sentences according to their relevance and then selecting a subset of these sentences that will form the summary. Several statistical methods can be applied in extractive summarization, and some of them are: term frequency-inverse document frequency (TF-IDF) and graph-

based algorithms such as TextRank. This approach is simple and often highly grammatically correct, but the summaries sometimes made are incoherent and without context. This is because the system extracts sentences without later modifying them to connect well together. It is useful for pieces of work where key information appears in distinct sections and research papers and reports are good examples of this type.

**3.2 Abstractive Summarization,** generates a new set of written sentences by understanding and paraphrasing the content contained in the original document. This method involves more sophistication because one requires the system to understand the meaning of the text and generate a summary cohesive and summarized. Mostly, abstractive summarization models rely on deep learning techniques, which are specialized to work with neural networks and Transformer models, especially BART-Bidirectional and Auto-Regressive Transformers and GPT-Generative Pre-trained Transformer. The primary reliance of these models is on the training that was given to them using large datasets of human-written summaries, allowing them to produce more human-like summaries compared to extractive methods. However, an issue with abstractive summarization is the problem of factual inaccuracy-for created text may contain information not contained in the source document.

**3.3 Transformer Models** have truly made document summarization really interesting by capturing long-range dependencies in the text, making them perfect for the summarization of long documents like PDFs and DOCX files. BART and T5, a Text-to-Text Transfer Transformer, showed amazing performance for abstractive summarization tasks. These models apply encoder-decoder architectures; the encoder computes the input document whereas the decoder generates the summary. It uses attention mechanisms to allow the model to pay attention to appropriate parts of the document, so a summary contains the most important information.

**3.4 Challenges in Summarization,** persist in the field of summarization, mainly in terms of

coherence, relevance, and fidelity when a document is summarized. The problem of ensuring that the selected sentences constitute coherent summaries afflicts extractive summarization. With abstractive methods, accuracy issues and hallucinations appear-hallucinations being the generation of information that is not present in the source document. Also, a large dataset and huge computational power would be required to train these models into effectiveness, limiting their applicability for most applications.

In terms of practical uses, automatic summarization is greatly helpful in areas where people are supposed to sit through a large amount of data and get the gist as soon as possible. For example, automatic summarization would be of great use in the field of law firms wherein lawyers need to go through vast contracts and case files to know the key points immediately. Similarly, summarization tools can be utilized by academic communities to summarize long research papers and immediately create abstracts for easy access with minimal effort. They're also being integrated into business workflows to summarize reports, meeting notes, and emails in general, therefore making for better productivity.

## 4. NAMED ENTITY RECOGNITION (NER)

Named Entity Recognition is part of the core NLP, but specifically deals with identifying key elements within a text, such as people names, organizations, locations, dates, and other specific entities. NER makes text otherwise unstructured more understandable and searchable in domains of interest, for example, legal documentation, academic research, healthcare, and business analytics. The structured information extracted from text using the technique of named entity recognition helps to assist multiple text-related tasks, such as information retrieval, document classification, and automated knowledge extraction.

**NER systems** have become remarkably different from the rule-based approaches to statistical models and now, most importantly, to the deep learning technique. Early versions of the NER model were greatly dependent on rule-based systems and hand-crafted features such as regular expressions and lexical patterns. Such systems required substantial domain expertise and were not all too easily adaptable to new

domains or even languages. There have been statistical models like Hidden Markov Models (HMMs) and Conditional Random Fields (CRFs) that addressed the limitations above. They are quite flexible in the sense that they learned patterns in the data using probabilistic methods rather than making simple determinations; however, dealing with complex ambiguous sentences was still problematic.

Performance on NER tasks has seen significant breakthroughs with the advent of deep learning models. Models based on neural networks especially require lesser human-designed features, which can automatically learn more complex data patterns. The family of Recurrent Neural Networks, particularly the Long Short-Term Memory, has enabled NER models to capture the sequential nature of text more effectively than ever before. However, the most critical advances were achieved through models of Transformer-based models, like BERT (Bidirectional Encoder Representations from Transformers) and RoBERTa (Robustly Optimized BERT).

- **BERT-based NER models** have been proven to be quite effective, in which they can, as a pre-trained model on large corpora of text, learn to understand the context in which an entity appears, more so when dealing with polysemous words and understanding the word dependencies in a sentence. Fine-tuning a pre-trained model in BERT-based systems on a specific NER dataset is actually used to obtain competitive state-of-the-art results regarding the accurate recognition of entities. Due to the self-attention mechanism of Transformer models, their ability to zoom in on the relevant parts of the sentence makes them very appropriate for NER tasks.

- **Applications of NER** are widespread across various domains. In the legal application, NER is used to identify names of parties, legal terms, dates, and case citations in enormous amounts of legal text for the simplification process of searching and review. It enables key concepts, researchers' names, and instituting through a sea of papers, hence easing better organization and retrieval in academic research. NER is applied in healthcare for extracting medical terms, patient names, and diagnoses from clinical notes and research papers.

One other major problem with NER is domain adaptation. Models pre-trained for NER might do well on general datasets but would likely struggle with texts that are particular to certain domains, like legal documents or medical reports, because of the difference in terminology and the structure of the texts from the structure of the training data. This can be improved by using fine-tuning on domain-specific datasets, but access to annotated data is often quite limiting in specialty domains. The issue here is also that of recognizing entities in low-resource languages where, by definition, large datasets for training are not available. This has been the focus of recent research: multilingual NER and cross-lingual transfer learning.

Named Entity Recognition is one of the very basic tasks in NLP, which was greatly enhanced with the advent of deep learning and the Transformer models. BERT and subsequent variants paved new benchmarks for the NER task in terms of achieving strong accuracy and contextual understanding. NER remains to be a vital tool for the structuring of unstructured data by applications in law, academia, healthcare, and so on through thousands of other domains, despite all the many current challenges in this research area such as domain adaptation, multilingual support, and ambiguity for entity. Future research will be on the performance of NER in specific domains and low-resource languages because the prospects and growth of this domain are likely to be anchored deeper as its applications grow by a more considerable extent across various industries.

## 5. ENGLISH TO HINDI LANGUAGE CONVERSION

Major strides have been taken with language conversion, also known as machine translation, using improved results from advancements in artificial intelligence and natural language processing. Converting from the English to Hindi language is one of the most critical endeavors in multilingual document processing, especially among countries like India, with Hindi being the most widely used language. The need to promote automatic translation systems is exigent in various areas like education, law, and government in the light of rising demands for cross-lingual communication. They should translate texts from one language into another with precision in meaning, context, and source document structure.

**5.1 Rule-based systems:** Which required hand-coding vast amounts of linguistic rules.

Such systems were rigid and not capable of withstanding the complexities and variations of natural languages, thus mostly producing translations that are grammatically correct but semantically wrong. Then came the SMT models that relied on huge corpora of parallel texts-paragraphs of sentences aligned in two languages to learn probabilistic word and phrase mappings. However, Hindi has much freer word order and a richer morphology than English, and SMT was severely at pain with these languages.

**5.2 Neural Machine Translation (NMT):** NMT models, particularly those with a deep learning basis, have transformed the paradigm landscape and offer more fluent as well as more contextually accurate translations. Continuous word or sentence and text representations learned by NMT models through neural networks enable them to understand more about the nuances of the language than earlier methods. The first successful NMT model was actually a Seq2Seq architecture with LSTM (Long Short-Term Memory) networks. This is an encoder-decoder framework where the source sentence, here being in English, is read by the encoder and compressed into a fixed-length vector, and this is used in the decoder to generate the target sentence in Hindi.

**5.3 Seq2Seq models : W**ere a step up, but still had its own limitations, particularly in the interpretation of long sentences; problems that would eventually be worked out with the advent of more advanced Transformer-based models, such as Google's BERT and OpenAI's GPT, which really brought translation quality to new heights. This model applies self-attention mechanisms, which enable the system to focus on the relevant parts of the sentence, thus making it reduce long-range dependencies and complex syntactic structures. It has been proven that transformer models are more effective for especially those languages in which word order is pretty flexible and meaning is heavily context-dependent, such as Hindi.

Problems in the translation of the **English text to Hindi** arise due to the drastically different grammars, syntax, and script systems used to write two such opposing languages. Hindi is an SOV, while the

English language is an SVO. Consequently, the translation model needs to alter the word order of the sentence and ensure that the meaning of the communication is somewhat similar. Hindi, in itself, has a developed morphology because the verb and nouns change according to the tense in which they have been pronounced. For example, the same word takes a different form if the person being referred to is male or female. Idioms and cultural references also raise many translation issues because literal translation seems to be a long way from what was intended.

The state-of-the-art models for translation from English to Hindi rely on Transformer architectures, such as Google's NMT model, the backbone of the Google Translate service. The **Google NMT model** relies on massive parallel corpora, allowing it to support learning over many language pairs, therefore enhancing the system's ability to generate grammatically correct and contextually relevant translations. Another open-source alternative that could be proposed for high-quality neural machine translation is **Hugging Face's MarianMT** models as well as **OpenNMT**. These models have been sensitized to complex grammatical rules as well as diverse linguistic contexts.

## 6. PLAGIARISM DETECTION TECHNIQUES

Detection of plagiarism is one of the most important areas in which information retrieval and natural language processing are applied. This is when instances of unacknowledged use or plagiarism of some other person's work, ideas, or intellectual property exist. Such detection and prevention have now become much needed not only in the science and technical fields but also in academic, journalism, and publishing environments where digital content keeps increasing. The primary function of plagiarism detection systems remains to safeguard the strength of academic performance, originality, and rights to intellectual property.

Plagiarism detection techniques can be broadly classified into three categories: **text-matching**, **semantic-based**, and **stylometric** methods.

> **6.1 Text matching techniques** represent one of the most common techniques used in plagiarism detection. These can be broadly defined as a comparison of the text to be

evaluated with a set of already existing documents for finding potential similarities. The simplest form of text matching is that of an exact string where the system searches for identical phrases or sentences that are available in the documents. More complex techniques depend on the tokenization and n-gram analysis, splitting up the text into smaller units (tokens or n-grams) and identify repeated subsequences. Methods like Levenshtein distance or Jaccard similarity may be utilized to define what are similarity measures between texts for the evaluation of possible plagiarism. Most applications taken into use today, such as Turnitin and Grammarly, function based on this approach, and report matches and direct links to real sources**.**

**6.2 Semantic-based techniques:** These exploit NLP and machine learning by attempting to determine the meaning of the text rather than its structure. In using these, such systems are enabled to identify plagiarism regardless of the fact that the wording was drastically different. For instance, Word Embeddings such as Word2Vec and GloVe, or Sentence Embeddings such as the Universal Sentence Encoder, may be leveraged to represent words or sentences in a continuous vector space. It captures semantic relations between words or sentences, so models can learn paraphrases, synonyms, and related concepts better, thus improving plagiarism detection.

**6.3 Stylistic methods** are another approach to the analysis of the style of a certain writer that might be searched for possible plagiarism cases. They depend on vocabulary richness, length of sentences, syntactic structure, and stylistic signs specific to the writers. As one of the stylistic profiling techniques, they may compare such a profile against established works by the author in order to determine the probability of the authorship. More usual than stylometric analysis are text-matching and semantic approaches, yet this can be useful in circumstances of much paraphrasing or modification to the text.

**-Outstanding challenges of Plagiarism Detection** in the treatment of multilingual texts as well as informality typical of other types of online content. Languages and idiomatic expressions show different syntactic structures, and thus direct comparisons will be difficult. In addition, collaborative writing tools and user-generated content led to more complex forms of plagiarism, such as patchwriting-the process by which authors compose texts that borrow from multiple sources but avoid direct quotation by using summaries, paraphrasings, or new words. Detection of this phenomenon requires advance techniques that could find the underlying structure and semantics of the text.

Deep learning techniques, during the last few years, have started to play a much more important role in plagiarism detection than before. Models of this type, based on the Transformer architecture such as BERT, have been well applied on tasks based on semantic similarity. In this work, by training those models on plagiarism detection datasets, we can substantially improve the performance enabling better detection of both copies and paraphrased content.

# 7. APPLICATIONS IN EDUCATION AND LAW

Document summarization and named entity recognition, language conversion technologies also have very important applications in education and law: whether it be plagiarism detection, technologies to summarize documents, named entity recognition, and language conversion make processes more efficient and set apart by their integrity for these sectors, providing vital support for the students, educators, legal professionals, and institutions.

## 7.1 Applications in Education

1. **Maintaining Academic Integrity**: Because plagiarism detection software like Turnitin and Grammarly is the most fundamental step in preserving the integrity of academic writing, checking students' submissions against an enormous body of already existing works might lead to the recognition of plagiarism situations thus encouraging original writings.

2. **Enhancing Learning Outcomes**: Thus, language translation technologies prove to be exceptionally useful in a multilingual education environment. Translation tools that take educational content from the mother tongue English and translate them in Hindi or any of the other regional languages make more students access and understand such learning content. That is most importantly so in a country with a diversification of linguistic backgrounds, as it increases the overall learning benefit in this regard.

3. **Streamlining Grading and Feedback**: Summarizing techniques for documents enable teachers to review large numbers of student work in a short time. By creating a summary of the particular assignment, a teacher can look at the vital points and arguments presented in the respective work so as to grade them more efficiently.

4. **Facilitating Research**: Named entity recognition helps students and researchers in the extraction of information pertaining to literature associated with their academic papers, articles, and online databases. NER technologies make the search process easy while researching, allowing students to easily find literature that can be used as supporting arguments for any query.

### 7.2 Applications in Law

1. **Legal Document Review**: There are broad specifications of documents that are supposed to be reviewed in the profession of the law. Plagiarism detection systems are helpful for a legal professional for their work being free from unintentionally duplicating pre-existent legal documents or precedents and thus by maintaining their integrity in writing and avoiding liabilities.

2. **Contract Analysis and Compliance**: Document summarization techniques assist attorneys in interpreting a contract, as well as legal document analysis. All summarization tools extract the key terms and clauses so that legal professionals are able to scan agreements more efficiently.

This is particularly crucial with regard to contract negotiation and compliance because it mandates careful attention to fine differences in legal language.

3. **Research and Case Law Retrieval**: The legal research process applies several applications, of which one is that of named entity recognition. NER systems automatically identify the relevant cases, statutes, and legal terms in the documents so that information is retrieved easily. It enables legal professionals to find precedents and references easily in order to assist in case preparation and formation of arguments.

4. **Language Translation in Legal Settings**: Language conversion technologies are largely instrumental in legal environments that use several languages. These tools enable legal documents and contracts as well as court proceedings to be translated to the tiniest detail, allowing all parties involved to get a very good understanding of the legal processes.

5. **Plagiarism Prevention in Legal Writing**: Originality reports of plagiarism detection systems can be availed by legal scholars and practitioners for checking the originality of their work. This is extremely relevant in publishing of academics, as law articles are expected to be sans uncredited sources. As a result, it ensures the credibility and reputation of the legal publications through originality.

## 8. CONCLUSION

The advancement of technology in the field of intelligent document processing has made tremendous changes in almost every sector-such as education and law-with the aid of tremendous automatic processes and algorithms that have identified ample opportunity areas in the world. Here, the four modules discussed are compression of PDF/DOCX files, NER, translation from English to Hindi, and plagiarism detection. With these modules, the researchers were able to observe the applications in the respective

domains. Here is the conclusion with the key findings and further research directions that could be followed for immediate implementation.

**Key Findings**

1. **Enhancing Academic Integrity**: Of utmost importance in all education institutions, plagiarism detection tools enable the avoidance of any infringement on ethical standards of writing. In addition to pointing out instances of unacknowledged use of other people's work, it also serves as an educative resource encouraging students to develop their own voice and ideas.

2. **Facilitating Research and Learning**: Many students and educators increasingly make use of named entity recognition and document summarization techniques in order to make research more streamlined. Such technologies automatically extract relevant information from the vast ocean of academic literature, thus allowing people to better traverse extended texts while they deepen their exposure to material.

3. **Bridging Language Barriers**: Language conversion technologies have proved worthy in multilingual environments. These can now offer students content and knowledge in native tongues; it really can make a difference, for instance, in very diverse regions where learners rely much on the understanding of what is being taught to them.

4. **Improving Legal Efficiency**: Legal-domain technologies, such as intelligent document processing, reduce the complexity in reviewing documents, contract analysis, and legal research. All these may leverage NLP techniques in order to automatically extract relevant information from complex legal texts that will assist lawyers and legal scholars with their work.

5. **Supporting Professional Development**: Advanced technologies for the processing of documents have many significant implications for professional development. They would most likely offer educators and other lawyers better tools that make their workflows easier without taking away the wasteful bureaucratic weight of offices, which
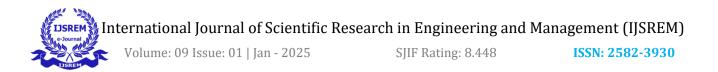
would leave more time for mentoring and advocacy.

## 9. FUTURE DIRECTIONS

Despite the progress made, several challenges remain that warrant further research and development:

1. **Domain-Specific Solutions**: Current technologies have proven generally effective, but domain-specific solutions are increasingly in demand, uniquely being able to satisfy the needs of specific fields of medicine, engineering, or legal expertise. Further development models that can account for unique jargon and nuances of these fields will enhance precision and applicability of such tools.

2. **Data Scarcity and Quality**: High-quality training and evaluation datasets in niche areas, such as named entity recognition, plagiarism detection, or other important tasks, will result in major challenges. Building broad open-access datasets on diverse languages, genres, and domains will be a collaborative effort toward further research in intelligent document processing.

3. **Ethical Considerations**: The rise of these technologies carries with it a host of ethical issues such as privacy, security of data, and algorithmic bias. To avoid situations where individual rights are trampled upon, or the use of these tools reinforces system inequalities, all these concerns need to be cautiously addressed.

4. **Integration and Interoperability**: Future work should be towards convergence of these document processing technologies into integrated systems that can work with existing platforms without a hitch. Such achievement of interoperability among different tools will promote better user experience and more coherent approaches toward document processing.

5. **User-Centric Design**: In terms of maximizing the impact, more user-centric design will happen. Thus, usability and accessibility will be considered in such designs. This means there will be much greater engagement among the targeted end-users in developing those tools and even in designing them properly to achieve several end-user needs and preferences by the students, the educators, and even the lawyers.

## REFERENCES

1] Maitri Patel and Dr Hemant D Vasava, " Natural Language Processing methods for Document Matching", International Journal for Modern Trends in Science and Technology, 2020, DOI:10.46501/IJMTST061271

2] Pradeepika Verma and Anshul Verma ,"A Review on Text Summarization Techniques" , 2020

3] Supriyono, Aji Prasetya Wibawa a, Suyono b, Fachrul Kurniawan, " A survey of text summarization: Techniques, evaluation and challenges" , 2024

4] Xu, Yiheng Li, Minghao Cui, Lei Huang, Shaohan ," LayoutLM: Pre-training of Text and Layout for Document Image Understanding", 2020

5] Asen Hikov ,Laura Murphy," Information retrieval from textual data: Harnessing large language models, retrieval augmented generation and prompt engineering", March 2024, DOI:10.69554/QAFE6376

6] Mohd Arsalan, " Transformers in Natural Language Processing: A Comprehensive Review", International Journal for Research in Applied Science and Engineering Technology, 2024.

7] Laxmi Reballiwar1, Sakshi Yergude , Vaidyavi, Sayli Birewar , Prof. Bhagyashree Karmarkar, "Language Translation Using Machine Learning", 2023

8] Ji Qi ,Wu Tongxin ,Yu Ting ,Dong Linxiao, " Power text information extraction based on multi-task learning", 2024, DOI:10.59782/sidr.v2i1.123

9] Ying Zhang ,Gang Xiao, " Named Entity Recognition Datasets: A Classification Framework", International Journal of Computational Intelligence Systems, 2024, DOI:10.1007/s44196-024-00456-1

10]C. Jiang, Y. Wang , J. Hu, J. Xu, " Power Entity Information Recognition Based on Deep Learning", 2021, DOI:10.13335/j.1000-3673.pst.2020.1678

11] Manish Balla, R Chowdary, "Generating Pre-Trained Transformers for Natural Language Processing ", 2023

12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," 2019.

13] Wolf, Thomas, Debut, Lysandre, Sanh, Victor et al ,"Transformers: State-of-the Art Natural Language Processing", 2020