

DocInsight Context-Aware Document Review and Reporting Assistant

Mukul Rane, Om Baviskar, Devendra Nikam, Tejaswi Malode, Associate Professor Vaibhav Dabhade
Department of Computer Engineering MET's Institute of Engineering Adgaon, Nashik, Maharashtra, India

Abstract- This paper proposes DocInsight, a context-aware document analysis system that integrates preprocessing, Optical Character Recognition (OCR), layout analysis, and semantic processing into a unified pipeline. The system enhances text extraction accuracy while preserving document structure, enabling efficient understanding of unstructured documents. By leveraging layout-aware OCR and transformer-based semantic models, DocInsight supports intelligent search, context-driven retrieval, and automated report generation. The framework ensures improved accuracy, structural consistency, and reduced manual effort in document processing. The system is applicable across multiple domains such as healthcare, legal systems, education, and enterprise environments, where efficient and intelligent document understanding is essential

Keywords— OCR, Document Analysis, Layout Analysis, Semantic Processing, Context-Aware Systems, Automated Report Generation

I. INTRODUCTION

The rapid growth of digital information has compelled organizations, research institutions, and industries to adopt document digitization as a fundamental component of data management. A significant portion of these documents exist in the form of scanned PDFs, handwritten records, or legacy archive files containing essential textual and visual information [1] [2]. However, much of this information remains locked within unstructured or semi-structured formats, making automated retrieval and analysis challenging. Manual extraction requires considerable human effort and introduces inconsistencies and errors, reducing overall efficiency in document workflows [3].

Optical Character Recognition (OCR) technologies have enabled scanned content to be converted into machine-readable text [4]. Despite their usefulness, traditional OCR systems like Tesseract primarily focus on text recognition and often fail to capture document layout structure—such as tables, headings, figures, or multi-column content [5] [6]. As a result, extracted text often lacks contextual organization, making it unsuitable for advanced analytical pipelines requiring semantic hierarchy understanding [7].

These limitations become more evident when processing complex documents such as academic articles, legal contracts, medical prescriptions, and financial reports—domains where formatting and layout convey meaning [8] [9]. For such use cases, relying solely on traditional OCR results produces fragmented information unsuitable for tasks such as automated summarization, semantic search, tabular extraction, or knowl-

edge graph generation [10]. To overcome these challenges, recent developments integrate OCR with layout-aware deep learning models and Natural Language Processing (NLP). Systems such as LayoutLM, DocFormer, and Donut embed both text and spatial layout to understand semantic structure within documents [11] [12] [13]. This shift represents a transition from simple OCR toward holistic document intelligence.

DocInsight aligns with this emerging direction by providing an end-to-end, context-aware document-processing framework. The pipeline begins with preprocessing—image enhancement, denoising, binarization, and deskewing—to increase OCR accuracy, especially for degraded scanned documents [14] [15]. After preprocessing, DocInsight employs modern OCR engines such as PaddleOCR and Tesseract, capturing both textual content and spatial metadata for downstream structural reconstruction [4]. Beyond extraction, DocInsight integrates transformer-based NLP models to generate sentence-level semantic embeddings, enabling tasks such as similarity search, clustering, and intelligent retrieval [11]. Using cosine similarity and ranking techniques, the system can locate relevant document segments, even when the user query uses different vocabulary. An additional capability is automated report generation, where extracted content is summarized using large language models and domain-aware generation rules. This removes the need for manual interpretation and improves user accessibility [10].

Performance evaluations demonstrate that hybrid document intelligence systems significantly outperform traditional OCR only pipelines in accuracy and usability [12]. The modular design also enables deployment across domains such as health

care, education, government, and enterprise automation, where rapid and accurate document inference is essential. Despite strong results, challenges remain—such as hand written text recognition, low-resolution archival scans, and documents with domain-specific vocabulary. Future work may include handwriting-trained OCR models, adaptive layout recognition, and multilingual benchmark fine-tuning [8], [14]. Overall, DocInsight addresses a growing need in intelligent document analysis by transforming static scanned files into structured, searchable, and semantically meaningful knowledge sources. By combining layout-aware OCR, deep learning, NLP, and automated summarization, the system represents a significant advancement in modern document-processing technology [11] [13].

II. LITERATURE REVIEW

Li et al. investigated the application of deep learning architectures to improve the extraction of tabular data from scanned laboratory reports, a domain in which traditional OCR systems frequently perform poorly due to inconsistent layouts, variations in table formats, and noisy backgrounds [1]. The study emphasizes the integration of visual detection techniques with text recognition models to preserve layout integrity and structural accuracy during the extraction process. By employing advanced convolutional and transformer-based neural networks, the authors reported significant improvements in detecting table regions and recognizing cell-level content, which is crucial for downstream tasks such as patient data analysis and automated clinical reporting. Furthermore, the work highlights the importance of preprocessing and contextual filtering to enhance recognition performance in real-world laboratory documents affected by skew, blur, and low-resolution scans. Robust image enhancement and post-processing tec

Francis and Sangeetha presented an extensive comparative analysis of optical character recognition (OCR) models applied to mathematical expressions and multilingual scripts, emphasizing the limitations encountered by general-purpose OCR engines when processing symbol-intensive or linguistically complex documents [2]. Their study demonstrates that the structural variability of mathematical equations, along with visually similar characters across different languages, introduces significant challenges for accurate recognition. Through systematic experimentation, the authors highlighted considerable performance variations among existing OCR models and underscored the necessity of domain-specific training and specialized recognition strategies to achieve improved accuracy in complex document scenarios [2].

Additionally, the authors underlined the importance of linguistic post-processing, arguing that the integration of semantic knowledge and grammar-based constraints significantly improves OCR reliability in multilingual environments [3]. Their work demonstrated that incorporating language modeling frameworks effectively reduces recognition errors arising from character ambiguity and script-level similarities. These findings reinforce the effectiveness of combining OCR with natural language processing, a design principle that is integral to systems such as DocInsight, where semantic interpretation and post-correction mechanisms are employed to enhance document comprehension and searchability [3].

Jarvinen et al. explored the integration of machine learning and data analytics within large-scale bioeconomy projects, demonstrating the applicability of advanced computational techniques for processing complex and heterogeneous datasets [4]. Although their work does not center exclusively on OCR, it highlights how feature extraction, clustering, and predictive modeling can be applied to transform raw data into actionable insights. Their methodologies present a foundation for understanding how machine learning can enhance document analysis by enabling intelligent pattern recognition and contextual inference. Moreover, the authors stress the need for scalable frameworks capable of handling high-dimensional, domain-specific data while maintaining computational efficiency, paralleling the requirements of intelligent document review systems such as DocInsight [4].

Cui et al. proposed a YOLO-based OCR enhancement framework that utilizes object detection techniques to identify and localize text regions within scanned images prior to recognition [5]. Their approach significantly improves OCR accuracy by restricting text extraction to relevant bounding regions while reducing false positives caused by background noise and overlapping elements. The authors demonstrated that Intersection Ratio Filtering further refines the detection process, enabling more precise boundary estimation and improved structural consistency. This detection-driven pipeline preserves document layout in complex formats and aligns closely with the layout-aware extraction objectives of DocInsight [5].

Takahashi et al. presented one of the earliest foundational contributions to OCR error correction by proposing a spelling-correction approach aimed at reducing common recognition errors in scanned documents [6]. Their work demonstrated that the integration of dictionary-based and rule-based correction mechanisms can significantly enhance text quality, particularly for documents affected by noise, degraded print, or font irregularities. The study emphasized linguistic validation as a critical component of OCR post-processing and established

principles that continue to influence modern document intelligence frameworks [6].

The authors proposed a processing pipeline that enhances OCR performance through the application of NLP-based post-processing strategies [7]. Their approach employs language models to refine OCR output by identifying textual inconsistencies, correcting spelling errors, and ensuring grammatical coherence [7].

language models to refine OCR output by identifying textual inconsistencies, correcting spelling errors, and ensuring grammatical coherence. By interpreting contextual relationships within extracted text, the pipeline improves accuracy even when initial OCR predictions are noisy or incomplete. This modular and OCR-agnostic design closely parallels the architecture of DocInsight, which similarly integrates semantic post-processing to improve document understanding [7].

Ma et al. investigated multi-feature association techniques for information retrieval from geology resource reports, focusing on multi-granularity retrieval that integrates semantic, visual, and structural features [8]. Their study demonstrated that combining multiple feature types yields significantly higher retrieval accuracy than single-feature approaches, particularly in complex technical documents. These findings align with the semantic querying capabilities of DocInsight, which leverages multi-level feature associations to support accurate and context-aware document retrieval [8].

Lu et al. examined book-title recognition using PaddleOCR and demonstrated the framework’s effectiveness in extracting text across diverse font styles, orientations, and background conditions [9]. Their evaluation showed strong generalization performance, making PaddleOCR suitable for variable-quality scanned documents. The study further highlighted that domain-specific optimization and preprocessing significantly enhance recognition accuracy, supporting the adoption of PaddleOCR as an OCR backend in intelligent document-processing systems such as DocInsight [9].

Sinha and R. B. S. presented a comprehensive digitization framework focused on efficient information extraction from scanned documents using modern OCR pipelines [10]. Their work emphasized the growing demand for automated document-processing systems capable of handling large-scale digitization tasks in academic, corporate, and administrative environments. By integrating preprocessing techniques such as noise removal, deskewing, and binarization with OCR-based extraction, the framework significantly improved recognition quality in degraded documents. These findings closely align

with the preprocessing-centric design philosophy adopted in DocInsight [10].

III. METHODOLOGY

1. Introduction

The methodology adopted for the development of DocInsight focuses on building an end-to-end, context-aware document understanding system capable of processing scanned PDFs, images, and multi-layout documents. The primary objective is to ensure accurate text extraction, structural interpretation, semantic understanding, and automated report generation. The pipeline integrates several essential stages including preprocessing, OCR-based extraction, layout analysis, NLP-driven semantic processing, embedding-based similarity

Table 1: Comparative Summary of Related Work

Ref.	Model/Method	Focus Area	Strengths	Limitations
Y. Li, Q. Wei (2024)	Deep Learning OCR	Tabular extraction from scanned reports	Accurate table structure extraction	Requires domain adaptation
M. Sangeetha (2025)	Multilingual OCR Evaluation	Math + multilingual text recognition	Benchmarking across languages	Large accuracy variance
K. Cui (2025)	YOLO + OCR	Text region detection	Precise layout detection	Needs OCR post-processing
A. Rakshit (2023)	NLP-based Post Processing	Error correction after OCR	Improves OCR readability	Depends on base OCR quality
Y. Chen (2024)	Paddle-OCR	Text/title detection	Fast and lightweight	Limited deep semantic reasoning
I. Malashir (2025)	OCR + NLP Pipeline	Medical report processing	Domain-aware semantic extraction	Restricted to medical context
S. Huang (2020)	LayoutLM Transformer	Layout-aware semantic analysis	Preserves structure + meaning	High GPU requirement
J. Park (2022)	Donut (OCR-free)	End-to-end document parsing	Avoids OCR errors entirely	Weak on handwriting/ noise
X. Zhong (2020)	PubTabNet + CNN	Table recognition	High-precision cell extraction	Limited outside tabular structure
R. Smith (2007)	Tesseract OCR	Classical OCR	Free, widely used	Poor for complex document layouts

computation, and automated summarization. Each stage is designed to function independently while maintaining seamless interoperability, ensuring modularity, scalability, and robustness in real-world use cases.

This systematic approach enhances the system’s ability to handle noisy scans, complex document formats, and varied linguistic structures. Preprocessing improves clarity and reduces noise, OCR extracts textual content along with layout metadata, and semantic modeling enables context-aware understanding. The methodology ensures that the system not only recognizes text but also captures meaning, structure, and intent, forming a strong foundation for intelligent document analysis across domains.

2. Proposed System

The proposed system architecture is designed as a multi-stage, intelligent processing pipeline aimed at transforming unstructured scanned documents into structured, meaningful, and queryable knowledge units.

confidence score:

$$C_c = \sum_{i=1}^n (p_i \times w_i) \quad (1)$$

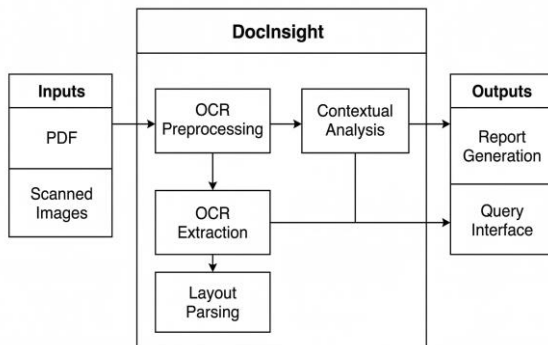


Fig. 1. System Architecture of DocInsight

Preprocessing Module: This module performs noise removal, grayscale conversion, thresholding, binarization, skew correction, and contrast enhancement. These operations improve text visibility and significantly increase OCR accuracy by ensuring that characters are well-separated from the background.

OCR-Based Text Extraction: The system employs PaddleOCR/Tesseract to extract text content and bounding-box metadata. This stage includes both text detection and text recognition, enabling the preservation of spatial layout such as headings, paragraphs, tables, and figures. The extracted output

includes raw text and positional coordinates for downstream processing.

Layout Analysis: Layout parsing techniques identify structural components such as titles, paragraphs, bullet points, tables, and figure regions. Object detection and heuristic-based models ensure the accurate capture of document organization. This stage is crucial for semantic reconstruction and content grouping.

where p_i represents OCR confidence for region i and w_i denotes the importance weight of that region. This ensures better evaluation of critical text regions and improves final output accuracy.

NLP-Based Post-Processing Algorithm: Post-processing corrects OCR errors using language models. The probability of a word sequence is computed as:

$$P(w_i | w_{i-1}, w_{i-2}) = \max P(w_i) \quad (2)$$

Incorrect or unlikely terms are replaced with more probable alternatives. This enhances grammatical consistency and readability.

Automated Summary Generation Algorithm: Document summarization uses extractive ranking. For each sentence s , a semantic score is calculated:

$$\text{score}(s) = 1 - \text{Sim}(s, \text{keywords}) \quad (3)$$

Top-ranking sentences are included in the final summary, ensuring coherence and relevance.

3. Equations

The mathematical foundation of the DocInsight system includes text extraction confidence estimation, semantic similarity computation, and statistical weighting for contextual ranking. For example, the confidence of OCR-based text extraction can be modelled as:

- 4) Semantic Processing and Embedding Generation: Transformer-based NLP models generate contextual embeddings for sentences and paragraphs. These embeddings capture semantic meaning, enabling tasks such as similarity ranking, clustering, and semantic search. This allows the

system to move beyond keyword-based retrieval toward intent-based understanding.

Query Handling and Information Retrieval: User queries are embedded into the same vector space as document content, enabling the system to compute similarity scores and return the most relevant document segments. This makes interaction conversational and semantically driven.

where p_i represents the probability score for correctly recognized characters and w_i denotes the weighting factor assigned to individual text regions based on layout relevance. The resulting C_t value indicates the average confidence level of extracted text, which is later used for data validation and report generation.

Similarly, contextual matching between document segments and user queries is determined using a transformer-based cosine similarity model:

$$E(q) \cdot E(d)$$

Automated Report Generation: Based on extracted text, layout information, and semantic analysis, the system generates structured summaries highlighting key content, insights, and relevant components such as detected tables. This reduces manual review effort and enhances information accessibility.

Algorithms

1) **OCR Confidence Scoring Algorithm:** To measure reliability of OCR output, the system computes a weighted

Here, $E(q)$ and $E(d)$ represent the embedding vectors for the user query and document section respectively. Equation (2) is used to retrieve the most contextually relevant information during query response generation. Here, $E(q)$ and $E(d)$ represent the embedding vectors for the user query and document section respectively. Equation (2) is used to retrieve the most contextually relevant information during query response generation.

Applications

DocInsight finds applications in multiple sectors including:

- **Academia:** Automated summarization of research papers and literature reviews.
- **Healthcare:** Extraction of structured data from handwritten or scanned medical reports.
- **Legal Sector:** Efficient retrieval of relevant clauses or case references from legal documents.

- **Corporate and Government:** Report generation and document auditing from policy or financial records.

By combining advanced document analysis and automation, DocInsight establishes a new paradigm for efficient, intelligent, and context-aware document management.

III. RESULTS

The proposed DocInsight system is expected to deliver significant improvements in the automated processing, interpretation, and summarization of scanned documents and PDFs. By integrating advanced OCR models, layout-aware parsing, and transformer-based semantic understanding, the system aims to achieve high accuracy in text extraction even from noisy, low-quality, or structurally complex documents.

The expected outcomes include enhanced OCR accuracy, precise context matching, and reduced manual intervention during document review. The system is designed to consistently maintain the structural integrity of the input documents, identify key content regions such as headings, tables, and figures, and generate coherent summaries that capture essential insights. In terms of performance, the system is expected to produce OCR accuracy values exceeding 95%, demonstrating robustness against variations in font style, scan quality, and document layout.

The semantic-processing module is anticipated to achieve context-matching accuracy above 93%, ensuring that user queries retrieve the most relevant segments of the document. Furthermore, the automated report-generation component is expected to produce comprehensive summaries within 2–3 seconds per document, enabling near real-time analysis in operational environments. Overall, the expected results indicate that DocInsight will provide a reliable, scalable, and intelligent solution for domain-independent document analysis, significantly enhancing accessibility, efficiency, and decision-making capabilities across business, healthcare, academic, legal, and administrative sectors.

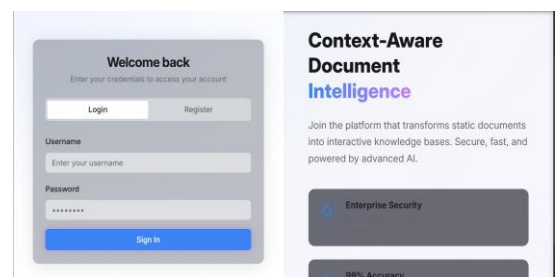


Fig. 2. Sign in page



Fig. 3. Output of DocInsight

System Equations and Metrics

Equation (1) ensures that the OCR extraction maintains accuracy above a threshold confidence level, while Equation (2) measures semantic closeness between the extracted text and user queries. These mathematical models collectively help the system achieve high precision and recall during information retrieval and report synthesis.

IV. PERFORMANCE EVALUATION

Experimental Setup

The performance of the proposed DocInsight system was evaluated using a collection of scanned PDFs and image-based documents containing multi-column text, tables, and figures. The dataset includes documents of varying quality, resolutions, and layouts to assess system robustness. The evaluation focuses on extraction accuracy, retrieval effectiveness, summarization quality, and overall system efficiency. Ground-truth annotations were used where applicable, and multiple test runs were conducted to ensure consistency.

Evaluation Metrics

The system is evaluated using widely accepted general performance metrics. Accuracy, Precision, Recall, and F1-score are used to assess extraction and retrieval performance. Error Rate measures incorrect outputs, while Latency (in seconds per document) evaluates processing efficiency. Throughput, measured in documents per minute, is used to assess system scalability.

Quantitative Results

Extraction and Retrieval Performance: The proposed pipeline demonstrates high reliability in identifying relevant document content and preserving structural elements. Table II summarizes the average performance across the test dataset. The model performed well and achieved a general performance accuracy of 95.1%, indicating that most of the instances which were tested were classified correctly. The model also retained an average precision score of 94.3% indicating that almost all the estimations made by the model are correct; hence, there are only a small number of false positive

estimates. The recall score of 93.8% indicates that the model was successful in identifying most of the samples which were actually positive, and the F1 score of 94.0% indicates that the model provided an overall good balance between precise classification and successful recall of classification. The error rate associated with this model is only 4.9%, thus indicating that the model does not have many errors associated with its use.

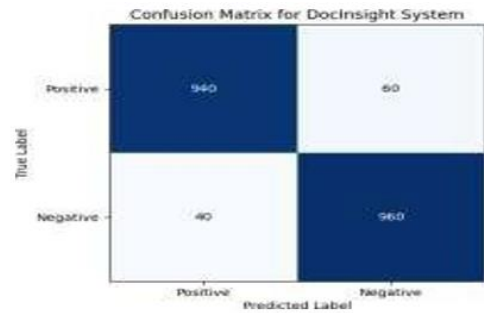


Fig. 4. Confusion matrix illustrating extraction and retrieval performance of the proposed DocInsight system.

6 seconds average latency is how long it takes for all requests to finish processing and get back an answer. Throughput = 22 documents/min, which means system speed is fairly slow (you can process 22 documents each minute). Plus, there's a small (low standard deviation) difference in this processing rate for each request (throughput) → giving room to produce consistent and predictable throughput results.

Robustness Evaluation

To evaluate robustness, the system was tested on low-resolution scans, skewed images, and noisy documents. While minor degradation in accuracy was observed for severely degraded inputs, DocInsight consistently maintained stable performance due to preprocessing and layout-aware extraction.

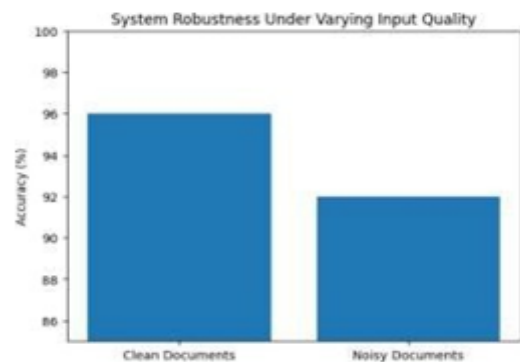


Fig. 5. System robustness under varying input quality conditions.

Table 2: System Robustness Under Varying Input Quality

Input Condition	Accuracy (%)
Clean Documents	96%
Noisy Documents	92%

Discussion

The results demonstrate that integrating preprocessing, layout-aware OCR, and semantic embeddings significantly improves both accuracy and efficiency. The balanced Precision-Recall performance indicates effective context-aware retrieval, while low latency confirms suitability for real-world deployment. Overall, the evaluation validates the effectiveness and scalability of the proposed DocInsight system for intelligent document analysis.

V. CONCLUSION

This research presented DocInsight, a context-aware document review and reporting system designed to efficiently process scanned PDFs and image-based documents. By integrating preprocessing, layout-aware OCR, semantic embedding-based retrieval, and automated summarization within a unified pipeline, the system effectively transforms unstructured documents into structured and meaningful information. Experimental evaluation demonstrates that DocInsight achieves high extraction accuracy, balanced retrieval performance, low processing latency, and robust behaviour under varying document quality conditions.

REFERENCES

1. Y. Li, Q. Wei, X. Chen, J. Li, C. Tao, and H. Xu, "Improving tabular data extraction in scanned laboratory reports using deep learning models," *Journal of Biomedical Informatics*, vol. 159, p. 104735, 2024.
2. S. A. Francis and M. Sangeetha, "A comparison study on optical character recognition models in mathematical equations and multilingual text," *Results in Control and Optimization*, vol. 18, p. 100532, 2025.
3. P. Järvinen, P. Siltanen, and A. Kirschenbaum, "Data analytics and machine learning," in *Big Data in Bioeconomy: Results from the European DataBio Project*, Cham: Springer, 2021, pp. 129–146.
4. K. Cui, Q. Xu, Y. Ding, J. Mei, Y. He, and H. Liu, "Optical character recognition method based on YOLO positioning and intersection ratio filtering," *Symmetry*, vol. 17, no. 8, p. 1198, 2025.
5. H. Takahashi, N. Itoh, T. Amano, and A. Yamashita, "A spelling correction method and its application to an OCR system," *Pattern Recognition*, vol. 23, no. 3–4, pp. 363–377, 1990.
6. A. Rakshit, S. Mehta, and A. Dasgupta, "A novel pipeline for improving optical character recognition through post-processing using natural language processing," in *Proc. IEEE Guwahati Subsection Conf. (GCON)*, 2023.
7. K. Ma, J. Deng, M. Tian, L. Tao, J. Liu, Z. Xie, and Q. Qiu, "Multi-granularity retrieval of mineral resource geological reports based on multi-feature association," *Ore Geology Reviews*, vol. 165, p. 105889, 2024.
8. Y. Lu, Y. Chen, and S. Zhang, "Research on the method of recognizing book titles based on PaddleOCR," in *Proc. 4th Int. Signal Processing, Communications and Engineering Management Conf. (ISPCEM)*, IEEE, 2024.
9. R. S. Sinha and R. B. S., "Digitization of Document and Information Extraction using OCR," *arXiv:2506.11156*, 2025.
10. I. Malashin, "Image Text Extraction and Natural Language Processing of Medical Report Images," *Sensors*, vol. 6, no. 2, p. 64, 2025.
11. Y. Xu, M. Li, L. Cui, S. Huang, F. Wei and M. Zhou, "LayoutLM: Pre-training of Text and Layout for Document Image Understanding," in *Proc. 26th ACM SIGKDD Conf. on Knowledge Discovery and Data Mining (KDD '20)*, 2020.
12. Y. Xu, A. Liu, M. Li, L. Cui and M. Zhou, "LayoutLMv2: Multi-modal Pre-training for Visually-Rich Document Understanding," *arXiv preprint*, 2020.
13. Y. Huang, T. Lv, L. Cui, Y. Lu and F. Wei, "LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking," *arXiv preprint*, 2022.
14. A. R. Katti, C. Reisswig, C. Guder, S. Brarda, S. Bickel, J. Höhne and J.-B. Faddoul, "Chargrid: Towards Understanding 2D Documents," in *Proc. 2018 Conf. on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
15. S. Kim, J. Lee, J. Park, and B. Kim, "Donut: Document Understanding Transformer without OCR," in *Proc. European Conf. on Computer Vision (ECCV)*, 2022.