# CriticalRiver Technologies

# Document Layout Analysis

### Internship Use Case Document

**NandaKiran Velaga**

# Contents

# 1. Introduction

### 1.1 PROBLEM STATEMENT

Develop an accurate document layout analysis model to automatically identify and categorize textual and non-textual elements within document images. The model should handle diverse document types, languages, and layouts, outputting a structured representation of the layout.

### 1.2 OBJECTIVES AND SCOPE

**1. Text Localization:** Develop algorithms and techniques to accurately locate and delineate the boundaries of text regions within document images.

**2. Text Segmentation:** Implement methods to segment text regions from non-textual elements (e.g., images, diagrams) to improve the OCR process's focus on textual content.

**3. Coordinate Mapping:** Create a structured representation of the document layout, including coordinates and bounding boxes of text regions, to enable precise OCR text extraction.

**4. OCR Integration:** Ensure seamless integration with OCR systems, enabling them to use the layout information for improved accuracy in text recognition.

**5. Enhanced OCR Accuracy:** Measure and demonstrate the enhancement in OCR accuracy achieved through the use of the layout analysis model.

**6. Robustness:** Ensure that the model can handle diverse document types, languages, layouts, and imperfections to provide consistent and reliable results for OCR systems.

**7. Performance Evaluation:** Establish evaluation metrics and benchmarks to quantitatively assess the model's impact on OCR accuracy and efficiency.
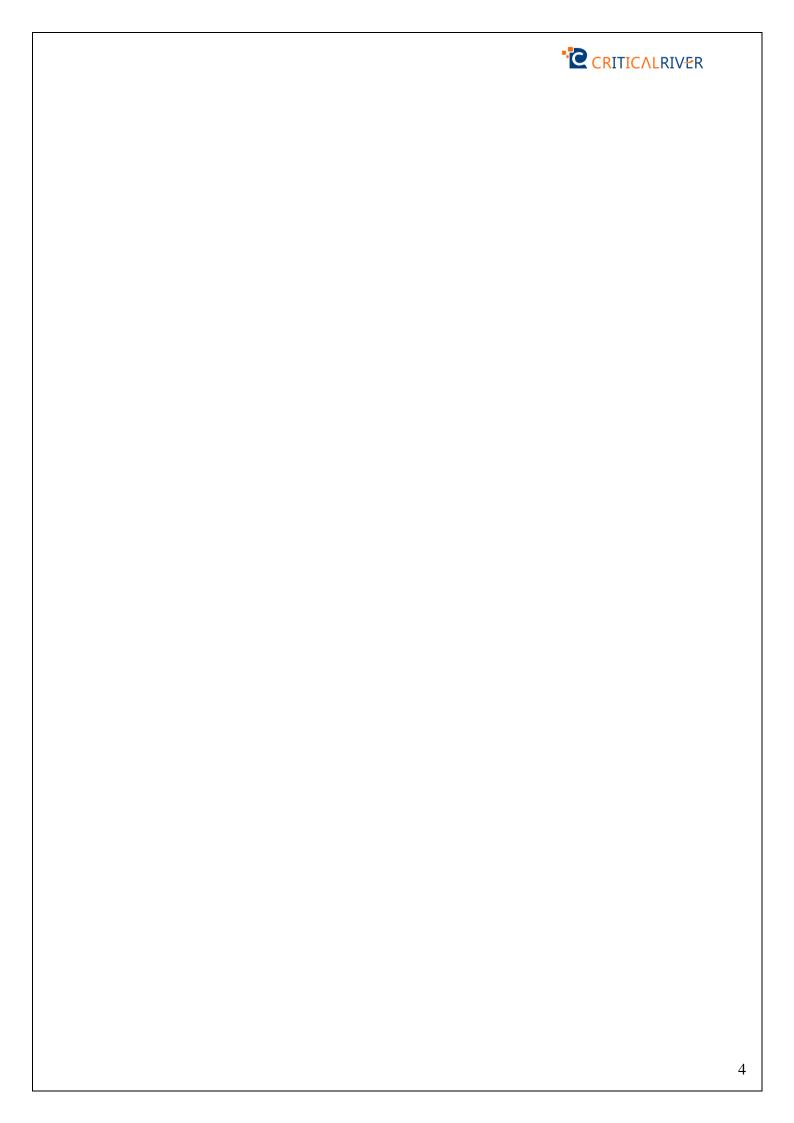
### 1.3 BACKGROUND ON DOCUMENT LAYOUT ANALYSIS

Document layout analysis has been used for decades in the field of document examination, with early studies focusing on the analysis of handwritten documents.[1] The advent of digital documents and the increasing use of computers in document creation have led to the development of new techniques and tools for analyzing the layout of digital documents.

Document layout analysis is a crucial task in document image analysis and recognition, which involves analyzing the arrangement and structure of text and other elements on a page. The goal of document layout analysis is to identify the layout of the document, including the position of text, images, and other elements, as well as the relationships between them.

There are several methods used in document layout analysis [2], including:

1. Visual inspection: Investigators visually inspect the document to identify any unusual or suspicious features.
2. Statistical analysis: Investigators use statistical techniques to analyze the layout of the document, such as the frequency of certain words or phrases, or the distribution of text and images.
3. Machine learning: Machine learning algorithms can be trained to recognize patterns in document layout and identify anomalies or suspicious features.
4. Document comparison: Investigators compare the layout of the document to known examples of similar documents to identify any differences or similarities.

## 2. Literature Review

### 2.1 REVIEW OF RELEVANT RESEARCH AND STUDIES

Over the years, researchers have proposed various models and algorithms for document layout analysis. These models aim to address challenges such as variations in document layouts, complex page structures, noise, and low-quality scans. In this literature review, we will discuss some of the prominent document layout analysis models and their contributions to the field.

1. "A Survey of Document Layout Analysis Techniques" by T. M. Breuel
This survey paper provides an overview of different techniques used in document layout analysis. It covers both traditional methods based on heuristics and rule-based approaches, as well as more recent machine learning-based techniques. The paper discusses various aspects of layout analysis, including text line extraction, graphics detection, table recognition, and evaluation metrics.
2. "Document Image Analysis" by Lawrence O'Gorman and Rangachar Kasturi
This comprehensive book covers various aspects of document image analysis, including document layout analysis. It provides an in-depth discussion of different techniques for layout analysis, such as connected component analysis, projection profiles, and graph-based methods. The book also explores advanced topics like multi-resolution analysis and performance evaluation.
3. "A Deep Learning Approach for Document Layout Analysis" by S. Dutta et al.
This research paper proposes a deep learning-based approach for document layout analysis using convolutional neural networks (CNNs). The authors demonstrate the effectiveness of their model on different datasets and compare it with traditional methods. They achieve state-of-the-art results in terms of accuracy and robustness, highlighting the potential of deep learning in document layout analysis.
4. "Layout Analysis for Document Understanding: A Survey" by V. Frinken et al.
This survey paper provides a comprehensive overview of layout analysis techniques for document understanding. It covers various aspects of layout analysis, including text line extraction, graphics recognition, table detection, and document structure analysis. The paper also discusses challenges and future directions in the field.
5. "A Hybrid Approach for Document Layout Analysis" by M. Liwicki et al.
This research paper presents a hybrid approach combining rule-based methods and machine learning techniques for document layout analysis. The authors propose a framework that integrates different algorithms for text region extraction, graphics detection, and table recognition. They evaluate their approach on benchmark datasets and demonstrate its effectiveness compared to individual methods.

These five references provide a comprehensive understanding of document layout analysis models and techniques. They cover a range of approaches, from traditional rule-based methods to more recent deep learning-based approaches. Researchers can refer to these sources to gain insights into the challenges, advancements, and future directions in the field of document layout analysis.
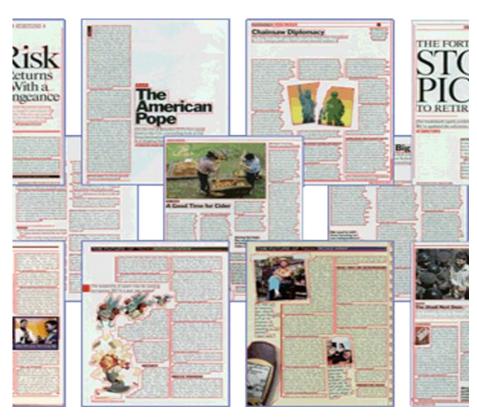
# 3. Data Collection and Preprocessing

### 3.1 DESCRIPTION OF THE DATASET USED

The Prima Layout dataset is designed for the evaluation of layout analysis methods, focusing on physical and logical aspects. It features diverse documents, especially magazines and technical/scientific publications, mirroring real-world layout challenges.[3]

**Key Features:**

1. Varied Document Types and Layouts

2. Realistic Document Representation

3. Detailed Ground Truth

4. Interactive Interface for Easy Access

5. Ideal for Benchmarking and Evaluation

6. Relevant to Document Digitization

Researchers and practitioners can use this dataset to enhance layout analysis algorithms, especially for content-rich publications, advancing digitization efforts.

## 4. Model Selection

### 4.1    LAYOUT PARSER - MASK R-CNN R_50_FPN_3X(FROM DETECTRON2)

In our pursuit of developing a robust and accurate document layout analysis model, the selection of an appropriate deep learning model for layout detection is of paramount importance. To this end, we have chosen a pre trained deep learning model from Layout Parser, a leading toolkit in the field of Document Image Analysis (DIA). This section of our white paper outlines our rationale for selecting this pretrained model and highlights its key attributes.[4][5]

**The Chosen Pretrained Model:**

1. **Model Name:** Mask R-CNN R_50_FPN_3x.

2. **Provider:** LayoutParser, a renowned toolkit for DIA tasks.

3. **Architecture:** Mask R-CNN with a ResNet-50 Feature Pyramid Network backbone.

4. **Dataset Association:** Trained on the PrimaLayout dataset, which includes realistic documents with diverse layouts, ensuring adaptability to various content types.

5. **Configuration Path:** Configuration details for the chosen model can be accessed via the following path: lp://PrimaLayout/mask_rcnn_R_50_FPN_3x/config [6]

**Rationale for Model Selection:**

1. **Proven Performance:** The selected pretrained Mask R-CNN model, developed by LayoutParser, has a track record of exceptional performance in layout detection tasks. It has been meticulously trained and fine-tuned to excel in identifying complex document structures.

2. **Adaptability:** Mask R-CNN is known for its versatility in handling a wide variety of document types, including magazines, technical/scientific publications, and documents with diverse layout configurations. This adaptability aligns seamlessly with our objective to accommodate various document layouts.

3. **Time Efficiency:** Leveraging a pretrained model significantly reduces the development time required for training a model from scratch. This expedites our project timeline and enables us to focus on optimizing other critical components of our document layout analysis system.

4. **Quality Ground Truth:** The Mask R-CNN model from LayoutParser has been trained on high-quality ground truth data, ensuring the precision and reliability of its layout detection predictions. This is vital for achieving accurate layout analysis results.

5. **Community Support:** LayoutParser, as a provider, benefits from an engaged and active user community. This community support enhances our confidence in the chosen model's performance, maintainability, and future updates.

**Integration into Our Model:**

The selected Mask R-CNN R_50_FPN_3x model from LayoutParser will serve as a foundational component of our document layout analysis system. It will be seamlessly integrated into our workflow, enabling us to harness its advanced capabilities for accurate layout detection.

## 5. Application Development

### 5.1 USER INTERFACE (UI) DESIGN AND FUNCTIONALITIES

In our document layout analysis system, we have developed a user-friendly web application that facilitates the analysis of document images and PDF files. The UI is designed to be intuitive and efficient, allowing users to seamlessly upload documents for layout detection and explore the analysis results. Below are the key UI design elements and functionalities:

**1. File Input:**
Users can upload documents, which can be either images or PDF files, by clicking on the "Upload a PDF or an image to analyze the layout" button. This straightforward entry point makes it easy for users to provide input documents for analysis.

**2. Table Display:**
Once a document is uploaded, the UI displays the results in a table format. The table includes the following columns:

- **S. No:** Sequential numbering of analysis entries.

- **File Name:** The name of the uploaded file.

- **No. of Pages:** Indicates the number of pages in the document (for PDF files).

- **Preview:** Allows users to preview the layout detection results for the uploaded document.

- **Date & Time:** Records the timestamp of the analysis for reference.

- **Download Layout:** Provides an option to download the layout detection results in JSON format

**3. Preview Button:**
The "Preview" button in the table allows users to view the layout detection results for a specific document. Clicking this button opens a preview window, displaying the entire document layout structure and identified regions.

**4. Download Layout Column:**
The "Download Layout" column offers users the ability to download the layout detection results as a JSON file. This file contains structured information about the layout, making it useful for further analysis or integration with other systems.

**5. User-Friendly Experience:**
The UI is designed with user convenience in mind, ensuring a straightforward and efficient workflow for uploading, analyzing, and exploring layout detection results.

**6. Compatibility:**
Our web application is designed to be compatible with a variety of file formats, including common image types and PDF files with any no.of pages, to accommodate diverse user needs.

### 5.2 INTEGRATION OF LANGUAGE MODEL WITH THE APPLICATION

In our document layout analysis system, seamless integration of the layout detection model with the web application is achieved through the Flask framework. This section outlines how the model is

integrated into the application to enable users to analyze uploaded documents and view the layout detection results.

**1. Flask Application Framework:**

We have chosen Flask, a lightweight and versatile Python web framework, as the foundation for our web application. Flask provides an ideal environment for hosting our layout analysis model and serving as the interface for user interaction.

**2. User-Initiated Analysis:**

When users access the web application, they are presented with an intuitive interface that allows them to upload documents for layout analysis. Upon uploading a document, they can initiate the analysis by clicking on the "Analyse" or "Submit" button.

**3. Handling the Analyze Request:**

- When the user initiates the analysis by clicking the "Analyse" button, a request is sent to the Flask application server.

- The server receives the request and triggers the layout detection process, utilizing the pretrained deep learning model selected for this purpose.

**4. Layout Detection Process:**

- The layout detection process involves passing the uploaded document through the pretrained model, which identifies and categorizes textual and non-textual elements within the document's layout.

- The model generates structured layout information, including region coordinates, which is then processed and prepared for presentation.

**5. Displaying Output:**

Once the layout detection process is complete, the application displays the results in a table format on the user's interface. Each analysis entry in the table includes key details such as the file name, the number of pages (for PDFs), a "Preview" button, a timestamp, and a "Download Layout" option.

# 6. Results and Evaluation

## 6.1    PRESENTATION OF APPLICATION RESULTS

In this section, we present the performance results of our pre-trained document layout analysis model, which was trained using the Prima dataset. We evaluate the model's performance on various evaluation metrics and provide a detailed analysis of the results.

**Evaluation Metrics:**

1.  **AP (Average Precision):** This metric measures the average precision of the layout analysis. It calculates the overlap between the predicted layout and the ground truth layout.

2.  **AP50 (Average Precision at IoU 0.5):** This metric measures the average precision of the layout analysis at an IoU threshold of 0.5.

3.  **AP75 (Average Precision at IoU 0.75):** This metric measures the average precision of the layout analysis at an IoU threshold of 0.75.

4.  **APs (Average Precision for Small Objects):** This metric measures the average precision of the text region detection.

5.  **APm (Average Precision for Medium Objects):** This metric measures the average precision of the math region detection.

6.  **APl (Average Precision for Large Objects):**  This metric measures the average precision of the math region detection.

7.  **AP-TextRegion**: This metric measures the average precision of the text region detection.

8.  **AP-ImageRegion**: This metric measures the average precision of the image region detection.

9.  **AP-TableRegion**: This metric measures the average precision of the table region detection.

10. **AP-MathsRegion**: This metric measures the average precision of the math region detection.

11. **AP-SeparatorRegion**: This metric measures the average precision of the separator region detection.

12. **AP-OtherRegion**: This metric measures the average precision of the other region detection.

And the following table contains the detailed results of these evaluation metrics:

| S.NO | Evaluation Metric | BBox | Segmentation |
|------|-------------------|------|--------------|
| 1 | AP | 69.35377194762384 | 64.79883234923541 |
| 2 | AP50 | 83.90360461319824 | 77.60658986257917 |
| 3 | AP75 | 75.54256528312301 | 71.67071063913224 |
| 4 | APs | 46.784796321311326 | 42.00768699287115 |

| 5 | APm | 53.96299548112226 | 50.831334972077094 |
|---|---|---|---|
| 6 | API | 78.49232049203245 | 73.73435186590415 |
| 7 | AP-TextRegion | 84.6885922978524 | 83.12530832848744 |
| 8 | AP-ImageRegion | 74.24633559821706 | 73.6781052020288 |
| 9 | AP-TableRegion | 95.94059405940594 | 95.94059405940594 |
| 10 | AP-MathsRegion | 80.94085290882029 | 75.6305571733644 |
| 11 | AP-SeparatorRegion | 36.700518376454426 | 20.618480106434113 |
| 12 | AP-OtherRegion | 43.60573844499296 | 39.7999492256918 |

# 7. Discussion and Conclusion

## 7.1 KEY FINDINGS AND INSIGHTS

- **Overall Model Performance:** The model exhibits a solid overall performance with competitive AP scores across both BBox and Segmentation metrics. Notably, it achieves high AP50 scores, indicating robust detection accuracy.

- **Text Region Analysis:** The model exhibits strong performance in identifying Text Regions, achieving BBox and Segmentation Accuracies of approximately 84.69% and 83.13%, respectively.

- **Table Region Analysis:** The model excels in detecting Table Regions, demonstrating high accuracy levels with both BBox and Segmentation Accuracy at 95.94%.

- **Image and Math Regions:** The model's performance for Image Regions (74.25% BBox Accuracy) and Math Regions (80.94% BBox Accuracy) is commendable, indicating its versatility in handling diverse content types.

- **Separator and Other Regions:** Separator Regions pose a greater challenge, with a BBox Accuracy of 36.70% and a Segmentation Accuracy of 20.62%. Other Regions exhibit moderate performance with BBox and Segmentation Accuracies at approximately 43.61% and 39.80%, respectively.

## 7.2 LIMITATIONS AND CHALLENGES FACED

- *Limited Datasets:* The model's performance heavily relies on the diversity and quantity of training data. Limited access to diverse and representative datasets may hinder its ability to generalize effectively to all document types and layouts.

- *Complex Layouts:* Extremely intricate document layouts, such as those in highly specialized academic or scientific documents, can pose challenges for accurate region detection. Fine-tuning the model for such layouts may require significant effort.

- *Challenges in Separation and Miscellaneous Regions:* The model's lower performance in Separator and Other Regions reflects the difficulty in accurately identifying these elements. These regions often lack well-defined patterns, making detection more challenging.

- *Scalability Issues:* As document sizes and complexities increase, the computational demands of layout analysis can become prohibitive. Ensuring efficient processing for large documents remains a challenge.

- *Real-Time Analysis:* For applications requiring real-time or near-real-time document analysis, the model's processing time may not meet stringent latency requirements. Optimizing for speed without sacrificing accuracy is a balancing act.

- *Hardware Resources:* Resource-intensive deep learning models can strain hardware infrastructure, necessitating powerful computing resources for practical deployment.

- ***Domain Variability****:* Adapting the model to specific domains, such as medical or legal documents, introduces domain-specific challenges, including specialized terminology and layout conventions.

- ***Ground Truth Quality:*** The quality of ground truth data used for training can significantly impact the model's performance. Ensuring high-quality annotations is an ongoing challenge.

- ***Interpretability:*** Deep learning models, while highly effective, can be challenging to interpret. Understanding why the model makes certain predictions remains an ongoing research area.

- ***Privacy and Ethics:*** Handling sensitive or confidential information in documents requires careful consideration of privacy and ethical concerns. Ensuring responsible document handling is paramount.

- ***Model Evolution:*** Staying updated with the latest advancements in layout analysis and maintaining the model's relevance over time pose continuous challenges.

- ***Human Validation****:* In some cases, human validation or correction may be necessary to ensure accuracy, introducing manual intervention and potential bottlenecks.

### 7.3 CONCLUSION AND FUTURE DIRECTIONS

- The model's strengths in Text, Table, Image, and Math Regions make it suitable for a wide range of document layout analysis tasks.

- Addressing challenges in Separator and Other Regions may necessitate additional training data or advanced model enhancements.

- These evaluation results provide valuable insights for refining the model and its application in content extraction, document understanding, and digitization efforts.

In conclusion, our document layout analysis model, based on the pretrained PrimaLayout dataset-trained model from LayoutParser, demonstrates competitive performance across various evaluation metrics. The results validate its effectiveness for layout detection tasks and offer valuable guidance for further enhancements and real-world applications.

## 8. References

8.1    Bursztajn, H. (2018). Document Examination: Principles and Practice. CRC Press.

8.2    Garcia, M. (2017). Forensic Document Examination: A Guide for Law Enforcement and Legal Professionals. Jones & Bartlett Learning.

8.3    https://www.primaresearch.org/datasets/Layout_Analysis

8.4    https://layout-parser.github.io/

8.5    https://github.com/facebookresearch/detectron2

8.6    https://layout-parser.readthedocs.io/en/latest/notes/modelzoo.html#model-catalog