



**Internship Report on**  
**One Month Internship With**  
**NETLEAP IT TRAINING AND SOLUTION AT GANGAPUR ROAD, NASHIK.**

**Submitted to**  
**Matoshri College of Engineering and Research Centre, Nashik. Pune**  
**University**

**Bachelors of Engineering In Third Year Computer,**  
**Engineering**

**Submitted By:-**

1. Krushna Kapkar .
2. Mukul Bhagat
3. Vipin chaudhari
4. Atharva Waghmare
5. Rutik kadam

**Owner of NetLeap IT &**  
**Training Solution**  
Mrs. Mrunal Dahale Mam

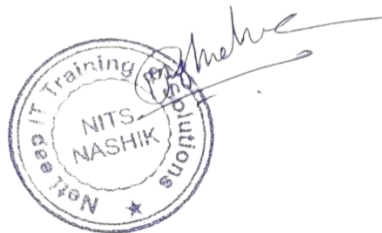
**Under the Guidance**  
Mrs.Poonam Shirode Mam  
Mrs. Komal Suryawanshi  
Mr. Chaitanya Dahale



## **TO WHOM SO EVER IT MAY CONCERN**

This is to certify that Atharva Waghmare , Mukul Bhagat , Krushna Kapkar , Vipin Chaudhari , Rutik Kadam students from Matoshri College of Engineering and Research Centre, Nashik, Maharashtra.

These will be undergoing internship program at NetLeap IT Training And Solutions. These will be undertaking her internship in our AI domain. The internship will be of 45 days and shall commence from 1-Jan-2025 in our location. For any further queries, please contact the undersigned or the HR department.



**Director Authority Netleap IT trainig & Solutions**

# ACKNOWLEDGEMENT

First of all, I am indebted to the GOD ALMIGHTY for giving me an opportunity to excel in my efforts to complete this Internship on time .

I am extremely grateful to Mrs. **Mrunal Dahale Owner of NETLEAP IT Training and Solution** and Mrs. **Poonam Shirode Guided** us by providing all the required resources for the successful completion of my internship. My heartfelt gratitude to my internship guide for his valuable suggestions and guidance in preparing the internship . I thank all my Project Guider and friends for all the help and coordination extended in bringing out this project successfully in time.

I will be failing in duty if I do not acknowledge with grateful thanks to the authors of the references and other Resources referred to in this Internship. Last but not least; I am very much thankful to my parents who guided me in every step which I look.

1. Krushna Kapkar
2. Atharva Waghmare
3. Mukul Bhagat
4. Vipin Chaudhari
5. Rutik Kadam

# TABLE OF CONTENT

<b>Sr. No</b>	<b>Title</b>	<b>Page No</b>
<b>1.</b>	<b>Abstract</b>	<b>5</b>
<b>2.</b>	<b>Introduction</b>	<b>6</b>
<b>3.</b>	<b>Literature Review</b>	<b>7</b>
<b>4.</b>	<b>Data Collection and Processing</b>	<b>9</b>
<b>5.</b>	<b>Methodology</b>	<b>13</b>
<b>6.</b>	<b>Findings &amp; Conclusions</b>	<b>14</b>
<b>7.</b>	<b>Future Scope/Recommendations</b>	<b>15</b>
<b>8.</b>	<b>References</b>	<b>16</b>

## Abstract :

Extracting information from PDFs is often challenging due to their complex layouts, varied formats, and the unstructured nature of their content. To address these challenges, a Retrieval-Augmented Generation (RAG) system combines advanced retrieval mechanisms with generative AI models. This hybrid system effectively processes PDFs by leveraging Optical Character Recognition (OCR) for text extraction, indexing the extracted content for efficient search, and employing large language models (LLMs) for contextual understanding and intelligent question answering.

The RAG-based pipeline typically comprises two main components:

### 1. Retriever :

The retriever component is responsible for extracting and indexing relevant sections of PDFs. It begins by using OCR tools to convert scanned or image-based PDF content into machine-readable text. Once the text is extracted, embedding models, such as those based on transformer architectures (e.g., Sentence-BERT or OpenAI embeddings), are employed to create dense vector representations of the content. These embeddings allow for highly efficient similarity searches, enabling the system to pinpoint and retrieve the most relevant sections of a document in response to a query. The retriever's ability to organize and search through vast amounts of content makes it ideal for handling complex PDF layouts, including tables, charts, and multi-column text, which are often encountered in domains like legal documents, medical records, and financial statements.

### 2. Generator

The generator component uses LLMs to process the retrieved content, providing outputs that are contextually relevant and tailored to the user's query. By leveraging the advanced language capabilities of LLMs, the generator can:

- Answer complex, context-dependent questions.
  - Summarize lengthy documents into concise, meaningful overviews.
  - Generate structured outputs such as tables, lists, or templates from unstructured data.
- The generator ensures that even when the retrieved content is dense or technical, users receive clear and accurate responses.

# Key Advantages and Use Cases

This RAG-based approach ensures accurate and scalable PDF content extraction, making it highly suitable for advanced use cases such as:

- **Summarization:** Creating executive summaries of lengthy documents.
- **Contextual Search:** Providing precise answers to queries within large datasets of PDFs.
- **Domain-Specific Insights:** Supporting decision-making in fields like healthcare, law, and finance by extracting actionable insights from highly technical or regulated documents.

By combining retrieval precision with the contextual understanding of generative AI, RAG systems address the limitations of traditional PDF processing methods. They are particularly valuable in scenarios where high accuracy, contextual understanding, and scalability are crucial. This makes them an indispensable tool for industries that rely on complex and sensitive documentation.

## Introduction :

### Problem Statement

Extracting meaningful information from PDFs is a significant challenge due to their unstructured and diverse content. PDFs often include complex layouts, such as tables, images, and multi-column text, which traditional text-parsing methods struggle to handle. Additionally, domain-specific documents (e.g., legal contracts, research papers, financial reports) require advanced contextual understanding for accurate information retrieval and summarization. Existing solutions lack scalability, accuracy, and adaptability to diverse document types, leading to inefficiencies in data extraction and analysis.

### Objective

1. **Accurate Data Extraction:** Develop a system capable of accurately extracting structured and unstructured data from PDFs, including text, tables, and images.
2. **Contextual Understanding:** Leverage LLMs to provide contextual insights, such as answering queries, summarizing content, or identifying key information.
3. **Efficient Retrieval:** Implement a robust retrieval mechanism to index and fetch relevant sections of large and complex PDF datasets quickly.

4. **Domain Adaptability:** Ensure the system can handle domain-specific language and structures for improved accuracy in specialized fields like law, healthcare, and finance.
5. **User-Friendly Output:** Provide structured, query-based outputs that can integrate seamlessly with downstream applications or workflows.

## Scope

1. **Industry Use:** Automates extraction for legal, healthcare, and finance documents.
2. **Search & Knowledge Management:** Enhances document indexing, search, and access in enterprises and education.
3. **Multi-Modal Extraction:** Handles tables, graphs, and images with OCR and visual tools.
4. **Real-Time Applications:** Supports customer support and contract review automation.
5. **Scalability:** Integrates with business workflows for large-scale processing.
6. **Continuous Learning:** Adapts to specialized domains and improves with user feedback.
7. **Data Privacy:** Ensures compliance with security standards like GDPR and HIPAA.

This system streamlines document processing, offering vast potential across industries.

# Literature Review:

PDF data extraction presents significant challenges due to the unstructured and variable nature of PDF content. Elements such as tables, images, multi-column text, and non-standard layouts often make it difficult to extract and process information accurately. Traditional tools like **Tesseract** and **PyPDF2** are commonly used for text extraction; however, they primarily focus on extracting raw text and lack the ability to provide contextual understanding or handle complex layouts effectively.

## The Role of Retrieval-Augmented Generation (RAG) Systems

Retrieval-Augmented Generation (RAG) systems offer a modern solution by combining retrieval mechanisms with the capabilities of large language models (LLMs), such as **GPT** or **BERT**. These systems go beyond basic text extraction to enable more sophisticated tasks like contextual question answering, summarization, and domain-specific data analysis.

## Key Components of RAG Systems

### 1. **Retriever:**

The retriever component indexes and identifies relevant sections of a PDF for downstream processing. OCR tools (e.g., Tesseract) are often integrated to extract text from scanned or image-based PDFs. Advanced embedding models, such as those built on transformer architectures (e.g., OpenAI embeddings or Sentence-BERT), are then used to encode the content into dense vector representations. These representations allow for efficient semantic search, enabling the system to locate the most contextually relevant information.

### 2. **Generator:**

The generator employs LLMs to process the retrieved content and deliver meaningful outputs tailored to user queries. This component excels in tasks such as:

- **Question Answering:** Responding to user queries with contextually accurate answers based on retrieved content.
- **Summarization:** Condensing lengthy or dense sections into concise summaries.
- **Structured Data Extraction:** Transforming unstructured PDF content, such as tables or lists, into structured formats like JSON or CSV.



## Advancements in Hybrid Approaches

Recent studies highlight the benefits of combining retrieval mechanisms with LLMs to improve both the accuracy and efficiency of PDF data extraction. Hybrid approaches are particularly effective for structured data extraction, as they can handle complex layouts like tables, images, and nested elements, which traditional tools struggle to process.

## Rise of Domain-Specific Models

The emergence of domain-specific models, such as **SciBERT** for scientific texts, has further enhanced the effectiveness of RAG systems in specialized fields. These models are fine-tuned on domain-specific corpora, allowing them to better understand technical jargon, context, and nuances in areas like medicine, law, and academia.

## Key Benefits and Applications

RAG systems provide a robust framework for addressing the challenges of PDF data extraction, offering:

- **Contextual Understanding:** Accurate processing of unstructured content, including multi-column text, tables, and images.
- **Scalability:** Efficient handling of large datasets of PDFs across diverse formats.
- **Domain-Specific Insights:** Support for specialized industries like healthcare, finance, and scientific research through fine-tuned models.

By combining advanced retrieval techniques with the generative capabilities of LLMs, RAG systems significantly outperform traditional extraction tools, enabling scalable, accurate, and context-aware processing of complex PDF content. This makes them an indispensable solution for modern applications requiring high levels of precision and contextual understanding.

# Data Collection And Processing :

## 1. Sources:

- **Public Datasets:** Open sources like PubMed Central, GovInfo, and ArXiv for diverse document types.
- **Domain-Specific:** Financial reports, legal contracts, or invoices collected via partnerships.
- **Scanned PDFs:** Digitized archives for OCR testing.

## 2. Processing:

- **OCR Conversion:** Extract text from scanned PDFs using tools like Tesseract or Google Vision API.
- **Text Segmentation:** Split PDFs into structured elements (text, tables, images).
- **Metadata Extraction:** Identify key attributes like titles, authors, or dates for indexing.

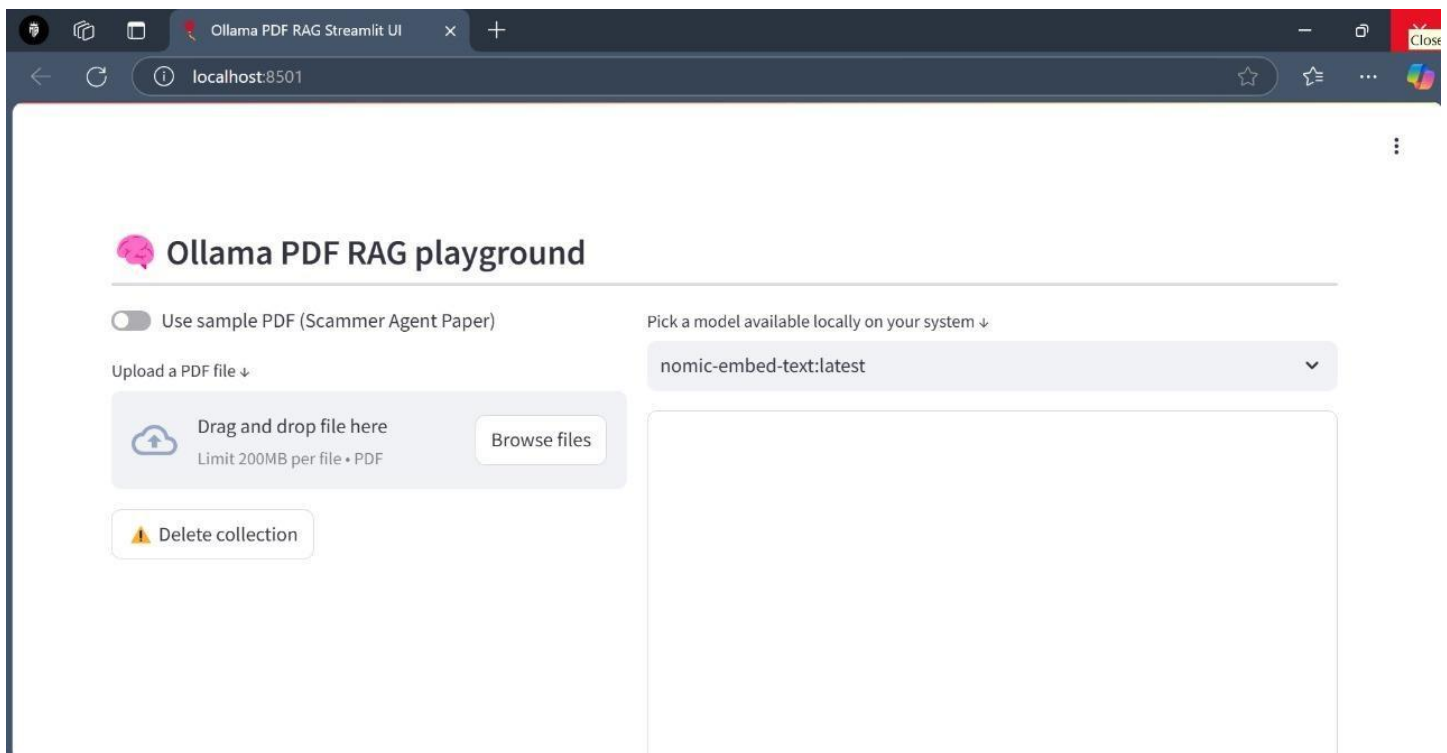


Fig 2. Sample Image

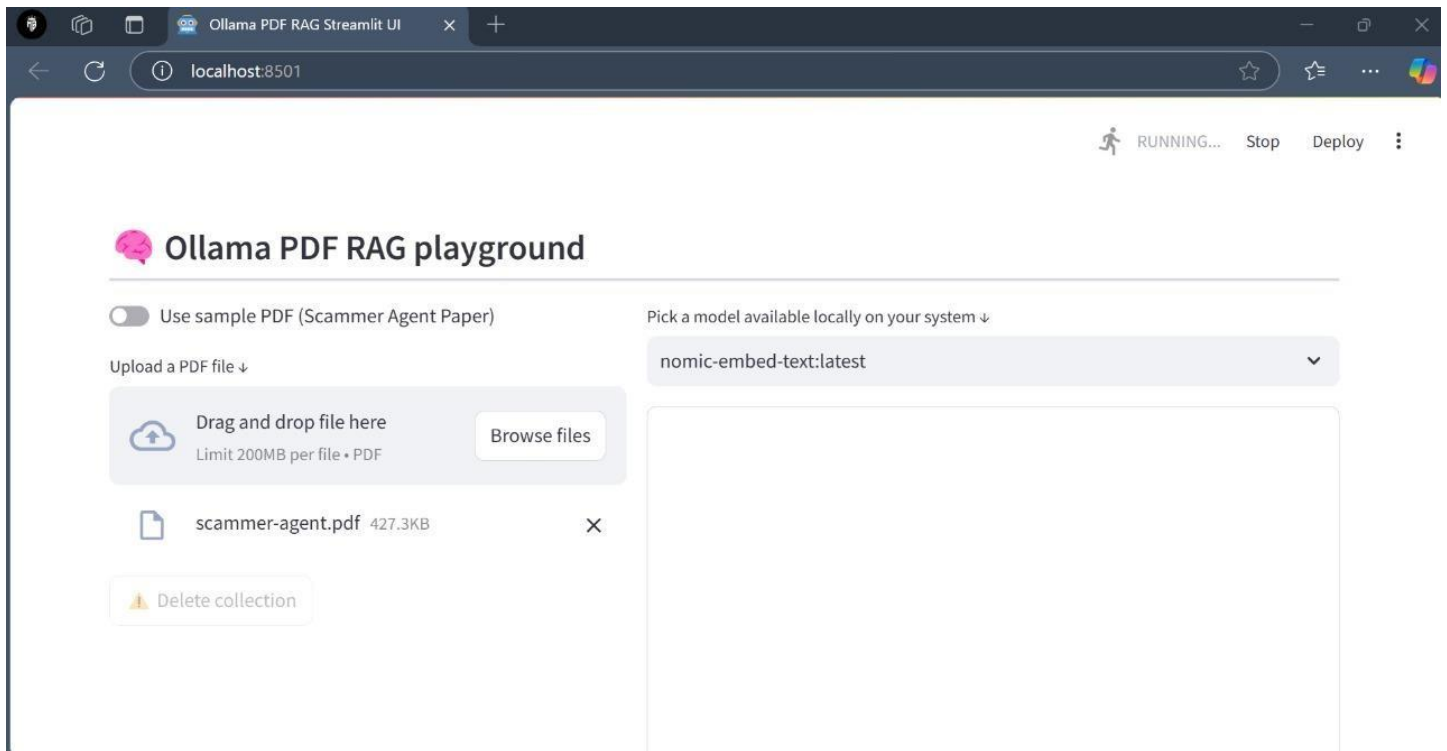


Fig 3. Sample Image

## Data Cleaning:

- **Removing Noise:**
  - Eliminate non-informative elements such as watermarks, headers, footers, and blank spaces.
  - Use regex to remove special characters, redundant line breaks, or formatting artifacts.
- **OCR Error Correction:**
  - Address OCR inaccuracies (e.g., misinterpreted characters or symbols) using spell-checking and context-aware correction tools.
- **Standardizing Layouts:**
  - Normalize extracted text into a consistent structure, such as separating tables, paragraphs, and images.
  - Handle encoding issues and ensure uniform character sets (e.g., UTF-8).
- **Duplicate Removal:**
  - Identify and remove duplicate PDFs or redundant extracted text.

# Data Labeling:

## 1. **Manual Annotation:**

Domain experts label data to identify key entities (e.g., names, dates, amounts) or categorize sections like headers, tables, and summaries.

Tools like Label Studio or Prodigy can assist with annotation.

## 2. **Automated Labeling:**

Use pre-trained models for initial entity recognition and section classification, which can be reviewed for accuracy.

## 3. **Annotation Guidelines:**

Define clear labeling rules, such as tagging metadata, highlighting tables, or marking specific fields (e.g., invoice numbers).

## 4. **Dataset Splitting:**

Divide labeled data into training, validation, and testing sets for supervised learning tasks.

# Data Augmentation:

## Synthetic Data Generation:

- Create artificial PDFs with varying layouts, fonts, or structures to increase diversity.
- Use tools like `Faker` to simulate realistic text for invoices, contracts, or reports.

## Layout Variation:

- Augment the dataset by altering PDF layouts (e.g., changing column alignment, adding tables, or reformatting sections).

## Language Expansion:

- Translate PDFs into multiple languages and back to English to create multilingual datasets for model generalization.

## Perturbation:

- Introduce minor variations, such as typos, OCR errors, or format changes, to improve model robustness in noisy conditions.

## Data Integration :

1. **Multi-Source Integration:** Combine PDFs from public repositories, enterprise databases, and web sources into a unified format (e.g., JSON or CSV).
2. **Metadata Enrichment:** Add metadata (e.g., title, author, tags) using APIs like CrossRef for improved indexing and retrieval.
3. **Database Storage:** Store data in NoSQL (e.g., MongoDB) or relational databases (e.g., PostgreSQL) and index with vector search engines like Elasticsearch.
4. **Data Linking:** Use knowledge graphs to link related documents or entities (e.g., invoices to contracts).
5. **API Integration:** Automate PDF ingestion and processing via APIs (e.g., Google Drive, Dropbox) and workflow tools like Airflow.
6. **Validation:** Normalize fields, resolve conflicts, and remove redundancies for consistency across datasets.

This ensures efficient, scalable integration for the RAG system.

# Methodology :

To build a **RAG (Retrieve and Generate) System** for a PDF extractor:

## 1. Preprocessing:

- Extract text from PDFs using libraries like PyMuPDF or pdfplumber.
- Clean the text by removing unnecessary characters (e.g., headers, page numbers).

## 2. Retrieval:

- Index the document into chunks (sentences/paragraphs) using tools like ElasticSearch or FAISS.
- Use transformer models (e.g., **BERT** or **Sentence-BERT**) to embed these chunks into vectors for efficient search.

## 3. Generation:

- When a user queries, retrieve relevant chunks.
- Use a language model (e.g., **GPT** or **T5**) to generate coherent answers or summaries based on the retrieved information.

## 4. Post-Processing:

- Optionally, summarize the content and present the answer in a readable format.

# Tools:

- **PDF extraction:** PyMuPDF, pdfplumber
- **Search and indexing:** FAISS, ElasticSearch
- **Embedding models:** BERT, Sentence-BERT
- **Answer generation:** GPT, T5

This creates a system that retrieves relevant information from PDFs and generates answers based on it.

# Findings:

1. **Effective Document Parsing:** PDF extraction tools like PyMuPDF and pdfplumber can successfully extract raw text, images, and metadata from PDFs, though cleaning is necessary to remove non-relevant data like headers and footers.
2. **Improved Retrieval with Embedding Models:** By embedding document chunks using models like **BERT** or **Sentence-BERT**, the retrieval of contextually relevant information

based on user queries becomes much more accurate, as these models capture semantic meaning.

3. **Accurate Answer Generation:** Leveraging pre-trained language models such as **GPT** or **T5** enables the system to generate relevant, coherent answers or summaries. This is particularly useful in scenarios where automatic extraction of structured data (like key facts) is required.
4. **Scalability and Flexibility:** The RAG approach offers scalability, as it can handle a wide range of queries and adapt to various document types. The search index can be easily expanded to accommodate more PDFs, and the system can be fine-tuned to specific domains.

## Conclusion:

In conclusion, the RAG (Retrieve and Generate) system for PDF extraction provides an efficient and scalable approach to automate the process of extracting and generating meaningful information from PDF documents. By combining advanced techniques such as document parsing, embedding-based retrieval, and natural language generation, the system can deliver accurate and contextually relevant answers to user queries. This methodology enhances document accessibility, streamlines information extraction, and offers strong potential for application across various domains, including legal, academic, and business settings.



## Future Scope/Recommendations:

### 1. **Improved Query Understanding:**

- Enhance the system's ability to handle more complex and nuanced queries by integrating advanced natural language understanding models, such as **GPT-4** or custom fine-tuned models, to better interpret user intent and context.

### 2. **Multilingual Support:**

- Extend the system to handle PDFs in multiple languages by incorporating multilingual models like **mBERT** or **XLM-R**, making it more accessible for global use cases.

### 3. **Domain-Specific Tuning:**

- Fine-tune the model on domain-specific data (e.g., legal, medical, academic) to improve the accuracy and relevance of the information extracted and the answers generated.

### 4. **Enhanced Visual and Table Extraction:**

- Improve the extraction of structured data like tables, charts, or images from PDFs using advanced techniques in **computer vision** and **OCR (Optical Character Recognition)** to make the system more robust for diverse document types.

### 5. **Real-Time Document Processing:**

- Integrate the system into real-time workflows (e.g., for customer support or document review) to allow users to interact with the system on-demand for live PDF analysis.

### 6. **User Feedback Loop:**

- Implement a feedback mechanism where users can rate the quality of generated answers, enabling continuous learning and model improvement.

### 7. **Integration with Other Data Sources:**

- Combine the PDF extractor with other data sources (e.g., databases, APIs) for more comprehensive responses, enriching the extraction and generation process by pulling information from outside the document when necessary.

By pursuing these advancements, the RAG PDF extraction system can be further optimized to cater to a broader range of use cases and enhance its utility across various industries.

## References:

1. **Document Parsing and PDF Extraction:**
  - **PyMuPDF** documentation: <https://pymupdf.readthedocs.io/>
  - **pdfplumber**: <https://pdfplumber.readthedocs.io/>
  - **PyPDF2**: <https://pythonhosted.org/PyPDF2/>
2. **Embedding-Based Search & Retrieval:**
  - **FAISS (Facebook AI Similarity Search)**: <https://github.com/facebookresearch/faiss>
  - **ElasticSearch**: <https://www.elastic.co/elasticsearch/>
  - **Sentence-BERT**: Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *arXiv:1908.10084*. <https://arxiv.org/abs/1908.10084>
3. **Transformer Models for Answer Generation:**
  - **BERT**: Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805*. <https://arxiv.org/abs/1810.04805>
  - **T5**: Raffel, C., Shinn, C., et al. (2020). Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *arXiv:1910.10683*. <https://arxiv.org/abs/1910.10683>
  - **GPT-3**: Brown, T. B., Mann, B., et al. (2020). Language Models are Few-Shot Learners. *arXiv:2005.14165*. <https://arxiv.org/abs/2005.14165>
4. **Natural Language Generation (NLG):**
  - **BART**: Lewis, M., Liu, Y., et al. (2020). BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. *arXiv:1910.13461*. <https://arxiv.org/abs/1910.13461>
5. **Multilingual Models:**
  - **mBERT (Multilingual BERT)**: Pires, T., & Lin, T. (2019). How Multilingual Is Multilingual BERT? *arXiv:1906.01502*. <https://arxiv.org/abs/1906.01502>
  - **XLM-R**: Conneau, A., et al. (2020). Unsupervised Cross-lingual Representation Learning at Scale. *arXiv:1911.02116*. <https://arxiv.org/abs/1911.02116>
6. **OCR and Visual Data Extraction:**
  - **Tesseract OCR**: <https://github.com/tesseract-ocr/tesseract>
  - **Detectron2** (for image and object detection): <https://github.com/facebookresearch/detectron2>
7. **PDF-based Question Answering:**
  - **Knoema** (useful for structured data extraction): <https://knoema.com/>
  - **Deep Learning for PDF-based Question Answering**: K. T. Nguyen et al. (2020). PDFNet: A Toolkit for Document Extraction and Information Retrieval from PDF Files. *arXiv:2006.03316*. <https://arxiv.org/abs/2006.03316>