# Project: DataRaptor

**DataRaptor** is a Docker-containerized web application designed for intelligent web scraping, data extraction, and analysis. It allows users to input a website URL, perform an automated dual scan, and catalog the extracted data into a PostgreSQL database. With integrated **AI-driven features** like **intelligent automation**, **data analysis**, **classification**, and **prediction**, DataRaptor transforms raw web data into actionable insights. Data can be exported in formats like JSON, XML, and RSS feeds, and users can schedule scraping tasks for ongoing updates.

## Technology Stack

1. **Programming Language**: **Python**
   - Leveraged for its rich ecosystem of web scraping libraries like **Selenium** and **BeautifulSoup**, as well as tools for AI and machine learning.
2. **Web Scraping Engine**: **Selenium**
   - Automates browser interactions, including dynamic content handling, user login, and website navigation, supporting complex, JavaScript-heavy websites.
3. **HTML Parser**: **BeautifulSoup**
   - Used for parsing and extracting structured and semi-structured data from HTML. It works in tandem with AI models to intelligently identify relevant information on web pages.
4. **Database**: **PostgreSQL**
   - Acts as the primary storage for all scraped and processed data. PostgreSQL ensures that data is securely stored and easily retrievable for querying and analysis.
5. **Containerization**: **Docker** and **docker-compose**
   - Ensures the application runs consistently across different environments. **docker-compose** orchestrates services like Selenium, PostgreSQL, and AI modules.
6. **AI Libraries**: **scikit-learn**, **TensorFlow**, **NLTK**
   - These libraries are used to integrate AI models for intelligent automation, data analysis, classification, and prediction.

## Key Features

1. **User Inputs URL & Credentials**:
   - Users provide the URL of the target website and any necessary login credentials.
2. **Dual Scanning Process**:
   - **First Scan**: Identifies login fields, enters credentials, and logs into the site.
   - **Second Scan**: After successful login, Selenium explores the site's structure, extracting relevant data from different sections like products, orders, and user details.
3. **AI-Driven Intelligent Automation**:

- ○ **Adaptive scraping**: AI models can detect changes in website structure (e.g., if the HTML changes) and adjust the scraping process dynamically.
- ○ **Predictive scraping**: AI can predict when the website's data will likely change and adjust scraping intervals accordingly, minimizing unnecessary scrapes.

4. **Data Classification and Categorization**:
   - ○ Scraped data is automatically classified and categorized using supervised machine learning models. For instance, products can be categorized into different types (e.g., electronics, clothing) based on their descriptions and metadata.

5. **Automated Data Analysis and Insights**:
   - ○ Machine learning models analyze scraped data to uncover trends, identify anomalies, and make predictions. For example, the system can predict price changes based on historical product data or flag unusual patterns in order histories.

6. **Handling Multiple Sections and Pagination**:
   - ○ Selenium automates navigation through different sections (e.g., products, orders) and handles pagination across multiple pages to ensure that all relevant data is collected.

7. **CAPTCHA Handling with 2Captcha and AI Models**:
   - ○ The app integrates **2Captcha** for solving CAPTCHA challenges during the scraping process. An optional AI-driven CAPTCHA solver can handle simpler CAPTCHA challenges locally, reducing reliance on third-party services.

8. **Data Storage in PostgreSQL**:
   - ○ All scraped data is securely stored in a PostgreSQL database. Users can query and export the data in various formats, such as JSON, XML, and CSV.

9. **RSS Feed Exposure**:
   - ○ The cataloged data can be exposed as an **RSS feed**, allowing users or systems to subscribe to updates and receive new data entries automatically. This makes the data accessible in real time for external systems.

10. **Scheduled Scraping**:
    - ○ Users can schedule scraping tasks to run automatically at set intervals. **Predictive scraping** powered by AI allows the system to adjust scraping frequency based on data patterns and predictions.

## AI-Powered Enhancements

1. **Intelligent Automation**:
   - ○ **Adaptive scraping** automatically adjusts to structural changes on the website using reinforcement learning models. The scraper can continue functioning even if the HTML structure of the website changes.
   - ○ **Error handling** is improved, as the system can dynamically respond to unexpected changes in website layout or content by re-training models in real-time.

2. **Data Classification**:
   - ○ Data is categorized based on patterns learned from previous examples. For example, product descriptions can be classified into predefined categories

(e.g., electronics, books, clothing), making it easier to organize and analyze large datasets.

3. **Data Analysis and Trend Prediction**:
    - AI models analyze trends within the data and make predictions. For instance, in an e-commerce scraping scenario, the system could track price trends over time and predict future price changes or stock availability.
    - **Anomaly detection** can flag unusual patterns in data (e.g., sudden spikes in pricing or outlier transactions) for further investigation.
4. **Predictive Scraping**:
    - The system uses historical data to predict the best times to scrape new information. For instance, if an e-commerce site typically updates its product listings every Monday, the system will prioritize scraping on that day.
5. **CAPTCHA Solving**:
    - AI-based CAPTCHA solvers (e.g., using **Optical Character Recognition (OCR)** for text-based CAPTCHA or **deep learning** for image-based CAPTCHA) can be integrated to handle challenges autonomously without relying on third-party services like 2Captcha.

## Workflow Overview

1. **Input URL and Credentials**:
    - Users input the website URL and login credentials (if necessary). AI models can assist in detecting form fields and inputs more intelligently if the website structure changes.
2. **Dual Scanning**:
    - The system performs a first scan for logging in and a second scan for collecting data. AI-driven automation helps adapt to changes in structure.
3. **Data Extraction and Classification**:
    - The extracted data is categorized using machine learning models, making it easier for users to organize and query specific datasets.
4. **CAPTCHA Handling**:
    - If CAPTCHA challenges are encountered, the system either uses **2Captcha** or locally implemented AI models to solve them automatically.
5. **Data Analysis and Prediction**:
    - AI models analyze the collected data, providing insights into trends, outliers, and making predictions for future behavior (e.g., price changes, user activity).
6. **Data Storage**:
    - The data is stored in **PostgreSQL** and made available in multiple formats (JSON, XML, CSV).
7. **RSS Feed Exposure**:
    - Users can subscribe to RSS feeds generated from the collected data to receive updates in real time.
8. **Scheduled Scraping with Predictive Insights**:
    - Scraping tasks can be scheduled, and AI models can predict when to scrape based on historical trends, optimizing the scraping process to ensure data is always up-to-date.

## Conclusion

**DataRaptor** is a robust, AI-enhanced web scraping and data analysis tool that goes beyond basic automation. It uses **intelligent automation** to adapt to changing websites, **machine learning** to classify and analyze data, and **predictive models** to anticipate data changes. By combining traditional web scraping techniques with advanced AI-driven features, DataRaptor ensures efficient, adaptable, and insightful data extraction for a wide range of use cases. Whether collecting data for business analysis or monitoring e-commerce trends, **DataRaptor** provides a comprehensive solution with real-time data exposure through **RSS feeds**.