# Di-Plast Matrix Data Extractor

Di-Plast Matrix Data Extractor (MDE) is a web-based application, which can be deployed on personal computer. It identifies document table regions on PDF documents using *Computer Vision based Deep Learning, especially Transfer Learning and Object Detection* algorithm. Then it extracts all textual data into text files by applying *Optical Character Recognition (OCR)* and also extracts tabular data separately in excel files using *Camelot* python package. It supports to transfer manufacturer names and corresponding technical datasheets names (or PDF filenames) to *MongoDB* database table for further processing.

The code can be downloaded from *GitHub* ( https://github.com/cslab-hub/MatrixDataExtractor ). **Only open-source software or library are used** in this application. MDE is primarily divided into 2 sections-

1. **Table Detection**: It provides code to train and test deep learning object detection model for document table detection task. The model is built mainly on *PyTorch Detectron2* library. The official support of Detectron2 library is available only on *Linux* OS (operating system). *Linux Desktop* version is recommended for normal users. Advanced users can use *Linux Server* version according to their choice.

2. **Backend**: It is the code of *Django* based web application, which provides a basic user interface to access the application functionalities for normal users.

First, build table detection model weight *(model_final.pth)*. Next, incorporate the model weight along with corresponding model description XML file

(*faster_rcnn_R_101_FPN_3x_config.yaml or uos_dip_config.yaml*) within Django application. The pre-requisite of MDE is given below-

1. For *normal users*, it is recommended to use *Linux Desktop* version (which comes with nice user interface), e.g. *Ubuntu OS*. The user interface is recommended for normal users
   - to create *Manufacturer* sub-folders
   - to store *Technical datasheets* (or PDF files) within Manufacturer sub-folders
   - check the model inference results, if table detection model correctly identifies document table regions on unseen document images or not.
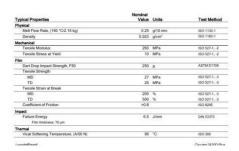
   *Advanced users* can use *Linux Server* version.

2. Please remember the difference between two terms-

   ***Document image***: Each PDF page converted into image format.



   ***Document Table Image***: Each document table on each document image.



3. Please check README file at GitHub page for installation. Install *MongoDB* and *Elastic Search* on Linux along with *Anaconda* environment and other libraries. An interface of MongoDB database (e.g. *MongoDB*

*Compass*) is recommended to access data from MongoDB. Elastic Search can be useful to search results based on textual query, which could be incorporated through future code development.

4. To train and test deep learning model, GPU enabled computer is recommended. Install *PyTorch 1.8.0 GPU* version and relevant *Detectron2* library for Table Detection section. For Backend section, you can use *PyTorch 1.8.0 CPU* version and relevant Detectron2 library.

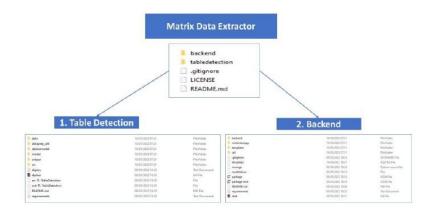The primary code structure is shown below-



Fig. Primary code structure

# Backend - Web Application

**Installation**:

Please check README file at GitHub page for installation. Optionally you can install *MongoDB Compass* tool.

**Folder Structure Overview**:

Normal users need to access *MatrixDataExtractor/backend/util* folder. Folder */util/prop* contains *MDE.xml* file, which is used for Django application configuration management. Folder */util/data/tabledet/modelweight* contains

- Faster R-CNN based object detection model description file: *faster_rcnn_R_101_FPN_3x_config.yaml* or *uos_dip_config.yaml*
- Table detection model weight: *model_final.pth*, which is built after table detection model training.

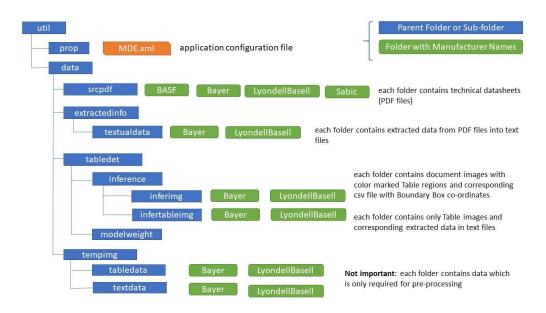The *MatrixDataExtractor/backend/util* folder structure is shown below-



Fig. Utility folder structure overview

**Pre-requisite for MDE web Application :**

1. Make sure */util/prop* folder contains *MDE.xml* file, which is used for MDE web application configuration.

2. Folder */util/data/tabledet/modelweight* contains

   - Model description file: *faster_rcnn_R_101_FPN_3x_config.yaml* or *uos_dip_config.yaml*

   - *Model weight file: model_final.pth*

These files are taken after Deep Learning Table Detection model training. Please refer *Model Training* sub-section of *Table Detection* section for more details.

**Run web application**:

1. Make sure *MongoDB and Elastic Search* services are installed on your Linux OS. You will find *start.sh* (shell) file in *MatrixDataExtractor/backend* folder. Execute the shell file with below command-

   **$ bash -i start.sh**

2. The web application starts on Anaconda environment *env_mde* by running the shell file. You can browse the web application by accessing **localhost:8000** url on your personal computer. If you want, you can give your preferred URL name at ALLOWED_HOSTS of *MatrixDataExtractor/backend/backend* folder's *setting.py* file.

**User guide of MDE web application**:

1. When the web application is running, you can see homepage as **Home** link (at left panel) along with other link descriptions. Several instructions are mentioned on webpage for simplicity.

2. Create sub-folders with *Manufacturer* names under */util/data/srcpdf* folder (s.g. BASF, Bayer, LyondellBasell). Keep corresponding PDF files within each folder. The sub-folders are created under */util/data/srcpdf* folder as below-
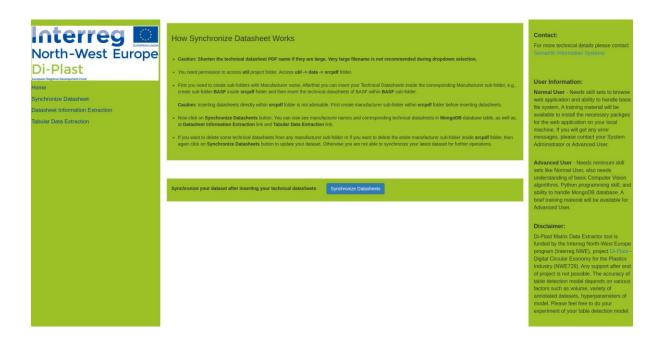


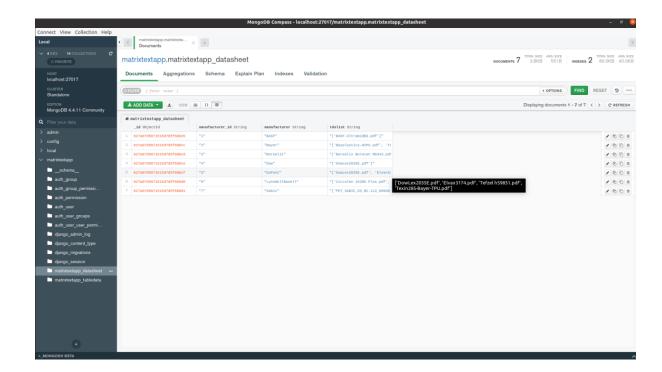The PDF files are kept within each sub-folder (e.g. LyondellBasell sub-folder) as below-



3. Go to **Synchronize Datasheet** link and click on **Synchronize Datasheets** button. It will synchronize all PDF files and corresponding sub-folders under */util/data/srcpdf* folder. It synchronizes
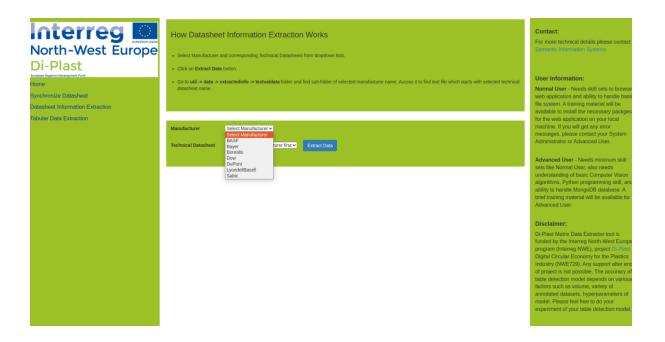
manufacturer names and corresponding PDF filenames in
***matrixtextapp_datasheet*** table in *MongoDB* database.
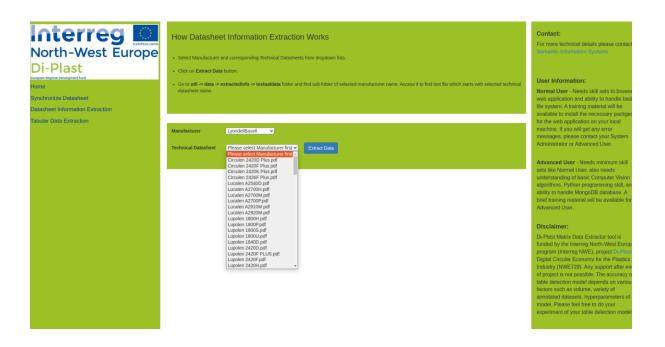


4. You can check manufacturer names and corresponding PDF filenames
   in ***matrixtextapp_datasheet*** table in *MongoDB* database. You can
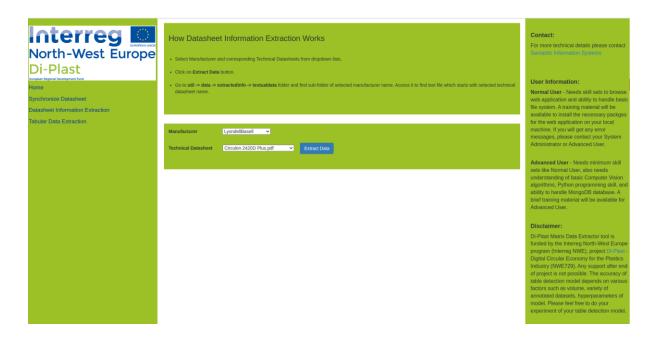   use *MongoDB Compass* tool to access the data.

5.  Go to **Datasheet Information Extraction** link to verify manufacturer names and corresponding technical datasheet names (or PDF filenames) available in dropdown menu. Select *Manufacturer* name first to access PDF files.
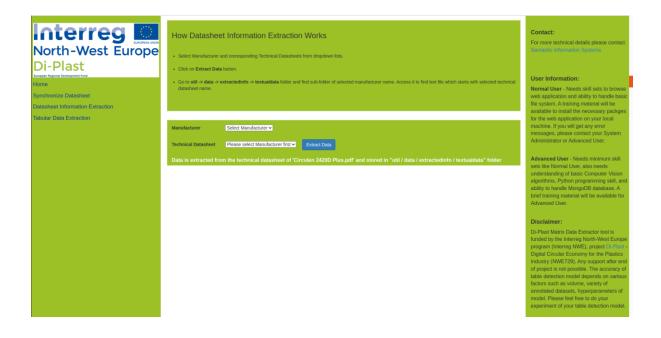


6.  Upon selection of *Manufacturer* name, you can get corresponding *Technical datasheet* names (or PDF filenames) in dropdown menu.

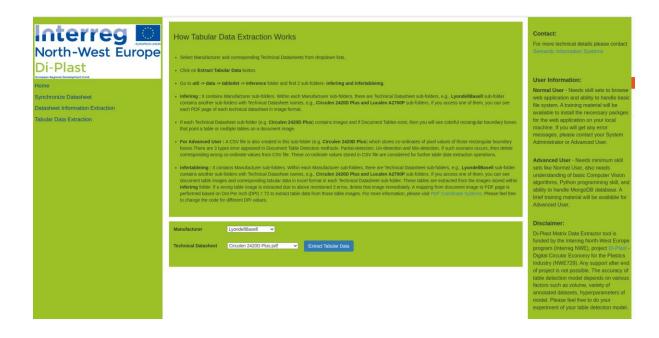7. Then select *Technical Datasheet* name (or PDF file) for further processing.



8. You can extract data from PDF files into textual format by clicking on **Extract Data** button. Data will be stored under */util/data/extractedinfo/textualdata* folder. You will get successful response (in white color) on webpage, if you have extracted data from PDF files as below-
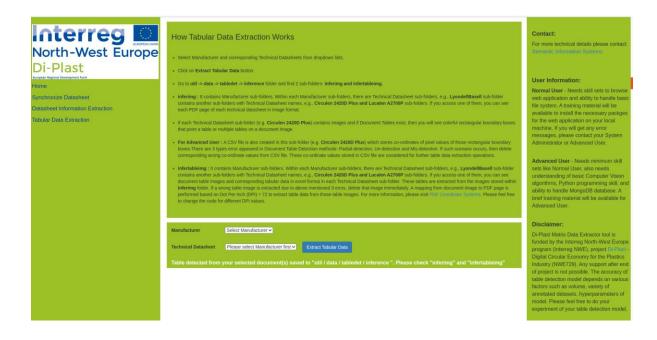
9. Please look above carefully that a sample notification is shown on web page after extracting 1 PDF document (e.g., Circulen 2420D Plus.pdf from LyondellBasell sub-folder is extracted) and stored PDF information in a text files.

10. Folder */util/data/extractedinfo/textualdata* contains *LyondellBasell* sub-folder, which also contains *Circulen 2420D Plus_preprocessed.txt* file as below. This unstructured textual information can be interesting for Elastic Search, Natural Language Processing (NLP) and Big Data technologies.



11. Go to **Tabular Data Extraction** link to click on **Extract Tabular Data** button.
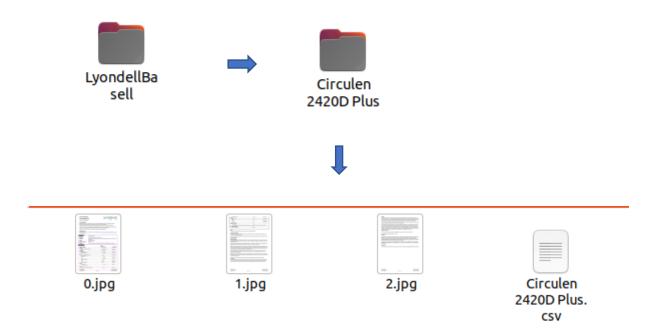
12. It identifies table regions of document images in a Rectangular Boundary Box (BBox) format under */util/data/tabledet/inference/inferimg* folder and stores BBox pixel information of document images in CSV files.

13. Simultaneously, you will also extract only document table images (cropped rectangular BBox region from document images) under /util/data/tabledet/inference/infertableimg folder and corresponding tabular data in excel files.

14. After clicking **Extract Tabular Data** button, you get successful response (in white color) on webpage. Now you access _document images and document table images both_ along with tabular data in excel files.
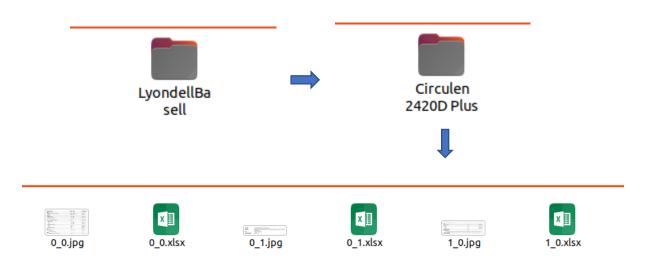


15. Above functionality involves _deep learning model inference_ to identify table regions in rectangular BBox format and cropped that regions to save table images. Folder /util/data/tabledet/inference/*inferimg* contains-

- _Manufacturer_ sub-folder

- *Technical Datasheet* sub-folder

- *Document images* along with BBox inference information in CSV file as below-



16. The document image pixels to PDF co-ordinates mapping is performed ( https://www.pdfscripting.com/public/PDF-Page-Coordinates.cfm ) to identify table regions on each PDF pages by considering DPI (dot per inch) value=72. DPI value is generally used to map digital images to physical pages (e.g. A4 page).

17. If you change DPI value other than 72, and your technical datasheets (or PDF files) are not A4 types, then feel free to adapt code changes to incorporate customized DPI value at *MatrixDataExtractor/backend/matrixtextapp/cv_basic_service.py*

18. When image pixel values to PDF co-ordinate values mapping is performed, *Camelot* python package is used with parameters *table_areas* and *flavor='stream'* to extract tabular data from PDF files in excel format.

19. Folder /util/data/tabledet/inference/*infertableimg* contains-

- *Manufacturer* sub-folder

- *Technical datasheet* sub-folder

- <u>*Document table images*</u> along with corresponding excel files as below-



20. Folder */util/data/tempimg* contains *tabledata* and *textdata* sub-folders for pre-processing purpose. Please delete all sub-folders and files within *tabledata* and *textdata* sub-folders when you finish your information extraction task.
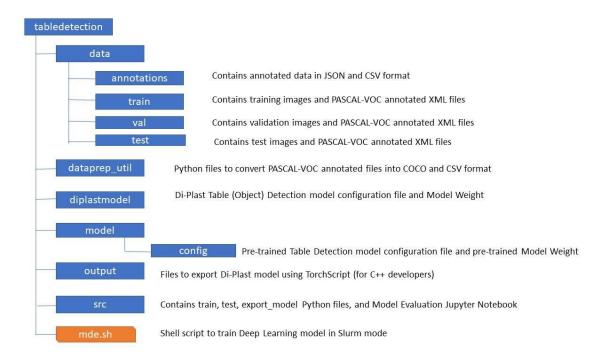
# Table Detection- Deep Learning Model

**Installation**:

Please check GitHub page for installation. The Anaconda environment ***env_mde*** needs to be created to train and to evaluate model (preferably in GPU server).

**Folder Structure**:

The folder structure is shown below-



**Configuration**:

1. Download pre-trained *TableBank* (*faster_rcnn_R_101_FPN_3x*) model config file and pre-trained model weight from *Layout-Parser* GitHub page (mentioned in *catalog.py* python file) from below URL-
   https://github.com/Layout-Parser/layout-parser/tree/main/src/layoutparser/models/detectron2

2. Save those files in *MatrixDataExtractor/tabledetection/model/configs* folder.

**Image Annotation**:

1. You can use any *Image Annotation* tool (e.g. *LabelImg*) to annotate **Table** images for Supervised Learning. If you have store data in PASCAL-VOC (XML) format, then you can convert XML annotated files into COCO (JSON) format. Also you need to convert annotated image information into CSV format for Di-Plast table detection model. The utility functions are referred in *tabledetection/dataprep_util* folder.

2. You can store JSON and CSV annotated information in *tabledetection/data/annotations* folder.

3. Keep your images and corresponding annotated XML files (PASCAL-VOC) in *tabledetection/data/train, tabledetection/data/val, tabledetection/data/test* folders.

**Model Training**:

1. In training, Di-Plast Table Detection model configuration file (*faster_rcnn_R_101_FPN_3x_config.yaml  or uos_dip_config.yaml*) and model training weight (*model_final.pth*) are saved in *tabledetection/diplastmodel* folder.

2. Run *tabledetection/src/train.py* script for model training. A shell script (*mde.sh*) is provided to train model in *Slurm* mode.

**Important Note**:

- You need to save *model_final.pth* and *faster_rcnn_R_101_FPN_3x_config.yaml  or uos_dip_config.yaml* files into */util/data/tabledet/ modelweight* folder in MDE web application for model inference.

**Model Evaluation**:

- Evaluate model by running *tabledetection/src/test.py* script. The Jupyter Notebook (*Eval_DiPlast_TableDetection_AP75 .ipynb*) is provided to evaluate the model and to visualize the inferred images.

## Optional- Export Model (for C++ Developers):

- *TorchScript*: Export model by running *tabledetection/src/export_model.py* script and saved model in *tabledetection/output* folder as *model.ts* format. This can be used in C++ development to infer Table images.

For more transfer learning based document table detection research work, please check below research paper-

Chowdhury, Arnab Ghosh, Nils Schut, and Martin Atzmüller. "A Hybrid Information Extraction Approach using Transfer Learning on Richly-Structured Documents." LWDA. 2021. ([http://ceur-ws.org/Vol-2993/paper-02.pdf](http://ceur-ws.org/Vol-2993/paper-02.pdf))