



TMA Solutions

# FINAL REPORT

NGUYEN DANG MAI THY

AI Center, TMA Innovations



# Introduction

---

I am Nguyen Dang Mai Thy, an intern from Batch 43 of the TMA Internship Program.

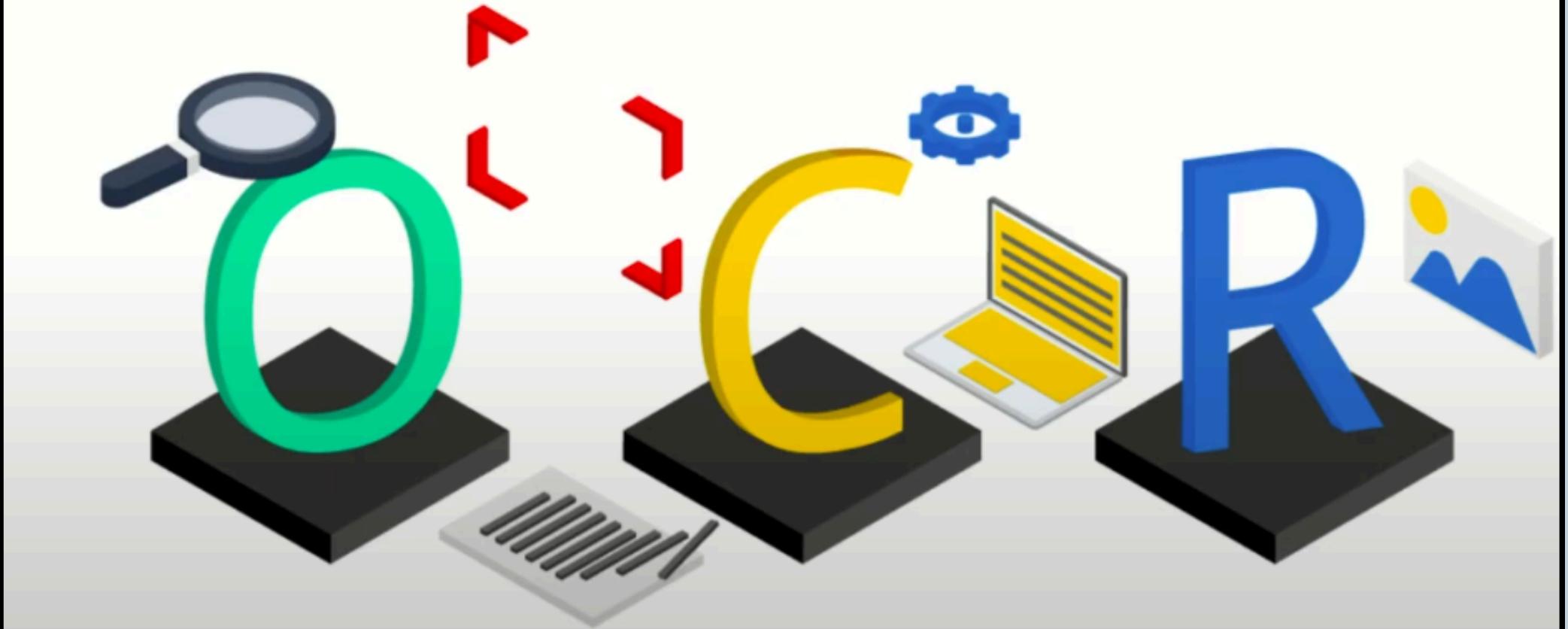
I was assigned to the AI Center, working in the Computer Vision domain, with Tran Trung Kien as my mentor and Nguyen Minh Nghia as the Project Manager.

I was allocated to the Dynamic Template OCR project and also provided support for several related projects.

## Welcome to OCR Data Extraction

This application allows you to create templates for OCR data extraction and extract information from images. Use the navigation sidebar to choose between creating a template and extracting information.

If you're new here, it's recommended to start with the guide to understand how to use the application effectively.



Dynamic Template OCR is a unique software that allows users to create custom OCR templates to process complex and unstructured documents.

Users can easily design and adjust the OCR templates to fit unique document types, such as contracts, reports, and invoices.

The custom OCR templates can identify and extract key data fields from documents, even when they lack clear structure.

# Requirements of Dynamic Template OCR

## Project Scope:

- The project covers processing of input images and PDF documents in Vietnamese and English.

## Project Requirements:

- Image preprocessing
- Document type recognition to classify into pre-existing templates
- Template creation including key-value pair extraction
- Table detection and table reconstruction
- Information extraction - Text OCR

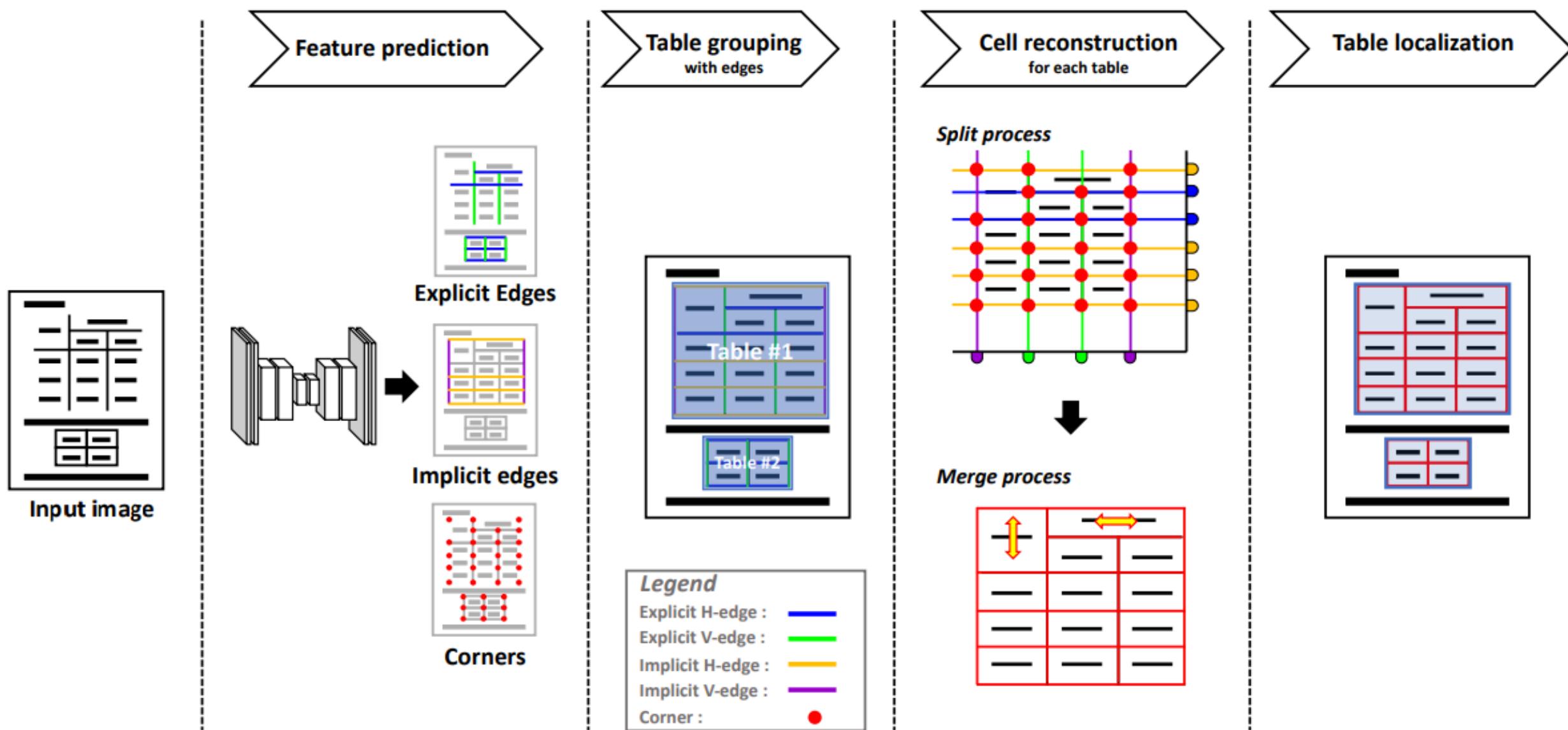
# Key deliverables

- Implement template grouping for template creation module
- Implement template grouping for information extraction
- Prepare data for table reconstruction
- Table Detection & Table Classification
- Table Reconstruction
- Integrate the project with Database and API
- Integrate the project with a user interface

# Result Achievements

- Research on Table Recognition (TR)
- Test & Compare TR repository
- Collect data to evaluate accuracy for TR model
- Draw Dynamic OCR full-flow
- Modify & Optimize TR code
- Collect data for training Label Medicine Box
- Research on PostgreSQL
- Research on FlaskAPI

# Research on Table Recognition



Document Understanding is a broad field that uses AI and automation to extract information from various document types, both structured and unstructured.

Table Recognition is a challenging problem within the broader field of Document Understanding. The goal is to accurately detect, recognize, and extract information from tabular data in documents, while also preserving the layout and structure of the table.

Key sub-tasks include Table Detection, Table Structure Recognition, Table Recognition, and the most complex, Table Understanding, which may involve extracting information from the table based on additional inputs like queries.

# Test &

Dataset	TD	TSR	TR	Number of sample	TSR Format
TableBank	✓	✓	✗	417K (TD), 144k (TSR)	Table Detection, Table2HTML
PubTabNet	✗	✓	✓	568K	Table2HTML
FinTabNet	✗	✓	✓	113K	Table2HTML
SciTSR	✗	✓	✓	15K	Text bounding box, Relative cell position, Cell Adjacency Relation
TIES_2.0	✗	✓	✓	Unbounded	Synthetic data, table-type classification, Table Structure to HTML
TabLex	✗	✓	✓	1M+	Structure information code, Content information code (similar to Table2HTML)
PubTable	✓	✓	✓	948K	Structure recognition, functional analysis, text content & bbox location
...					

# Compare TR repository

Dataset	TSR Format	Note
PubTabNet	Table2HTML	1/2 are span-cells + majority are borderless ⇒ when the input is full-border, the accuracy is low (majority of Vietnamese data is full-border)
TableBank	Table Detection, Table2HTML	
SciTSR	Text bounding box, Relative cell position, Cell Adjacency Relation	
FinbTabNet	Table2HTML	
TNCR	5 class	

For the pipeline, the Table Detection model can be trained on a composite dataset combining resources like TableBank, FinTabNet, PubLayNet, and TNCR. The model will detect a single table class, with a separate model used to classify the table type in a subsequent step.

# Test & Compare TR repository

Repository	Advantages	Disadvantages
PaddleOCR-ppstrcter	<ul style="list-style-type: none"> <li>- Lightweight and optimized model</li> <li>- Supports direction correction</li> </ul>	<ul style="list-style-type: none"> <li>- Limited dataset</li> <li>- Difficult to develop and expand</li> <li>- Challenging post-processing</li> <li>- Poor Vietnamese language support</li> </ul>
TableMaster, MTL TabNet, EDD-third-party, tsr-convstem	N/A	<ul style="list-style-type: none"> <li>- Difficult to create dataset</li> <li>- Heavy models</li> <li>- Extremely difficult and infeasible to create more data</li> </ul>
GraphTSR	N/A	<ul style="list-style-type: none"> <li>- Difficult to prepare the data</li> <li>- Unrealistic assumptions about the input</li> </ul>
deepdoctection	<ul style="list-style-type: none"> <li>- Can be developed and expanded the best as it allows customization of the format and use of custom data</li> </ul>	<ul style="list-style-type: none"> <li>- Default uses the PubTabNet dataset, so it cannot detect borderless tables, but can be retrained with custom data to improve</li> <li>- Works well on the demo app, good Vietnamese language support</li> </ul>
Davar-OCR	N/A	<ul style="list-style-type: none"> <li>- Quite complex design with many sub-modules</li> </ul>
DeepTSR, TableNet	N/A	<ul style="list-style-type: none"> <li>- The model is too simple, unable to detect span-cells</li> </ul>

The comparison indicates that while DeepDocDetection has strong table detection and cell detection, its table reconstruction is weaker and it cannot handle span cells.

After further research, TRACE appears to be the best performing in terms of detecting and redrawing span cells, as well as overall table reconstruction, outperforming all the other remaining repositories.



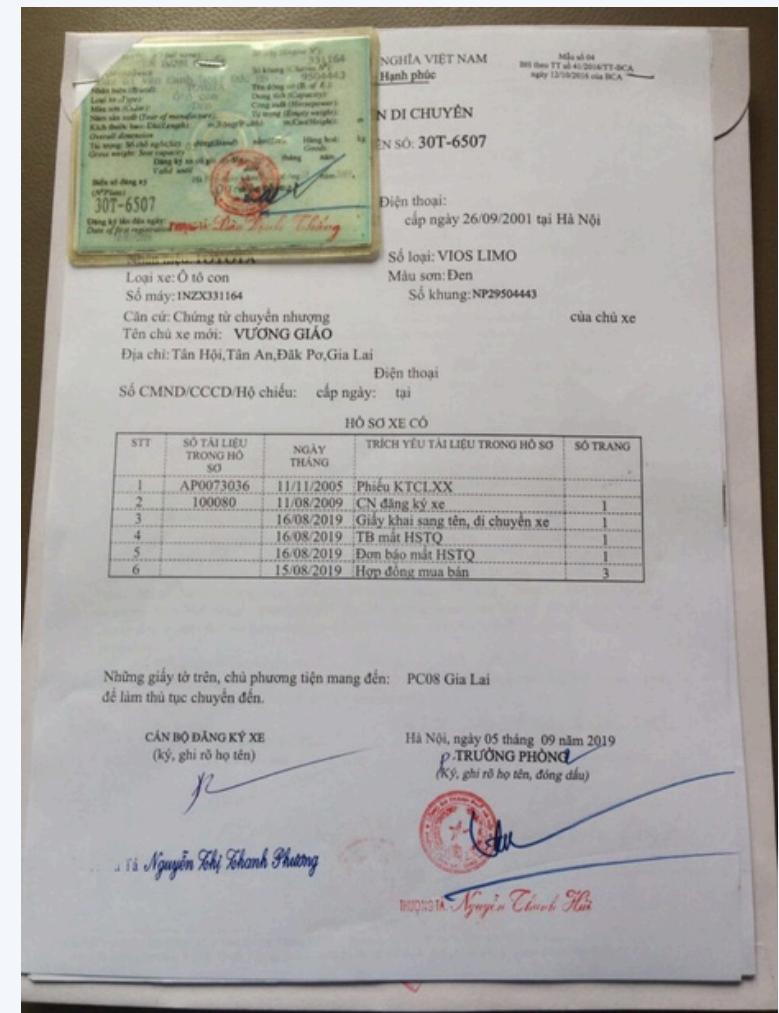
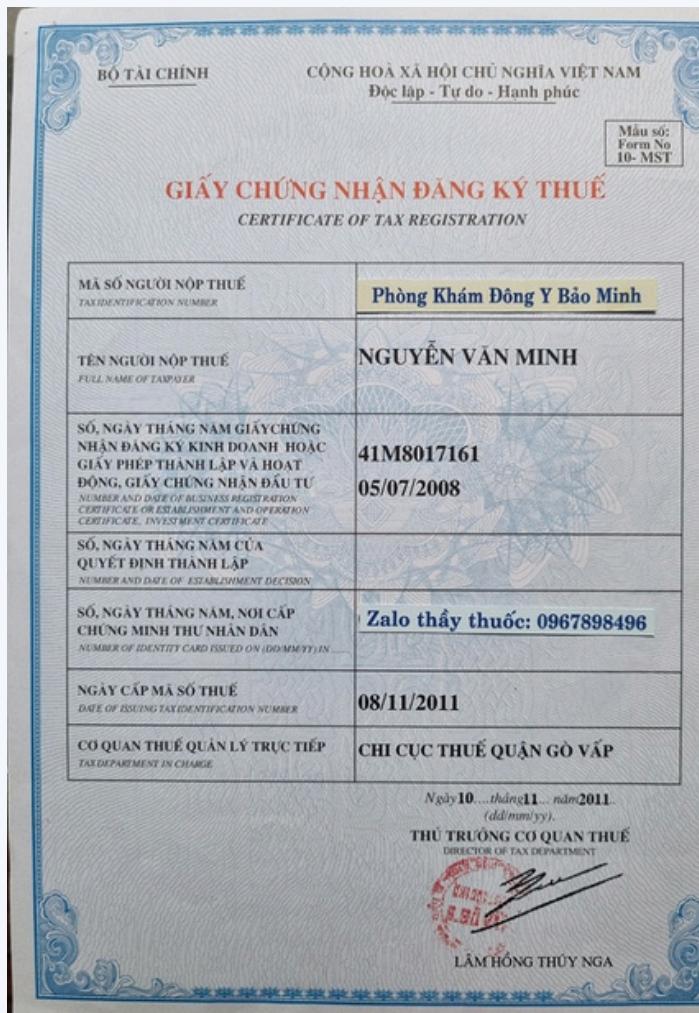
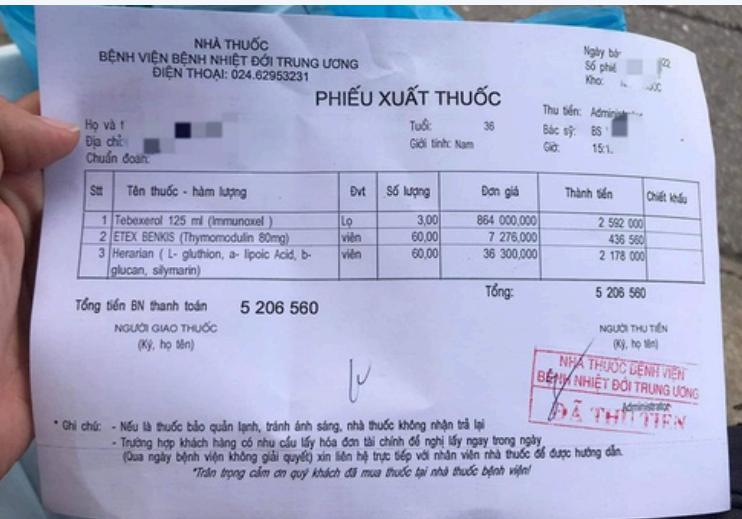
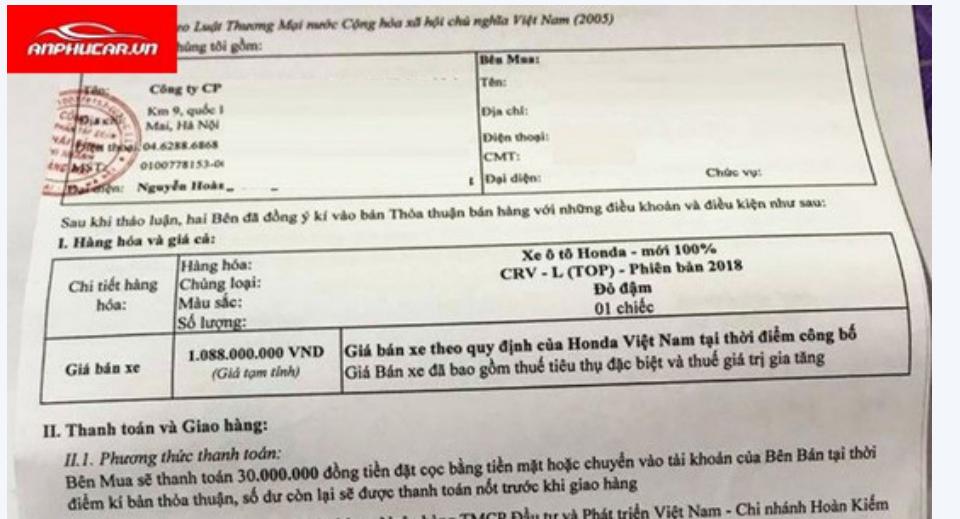
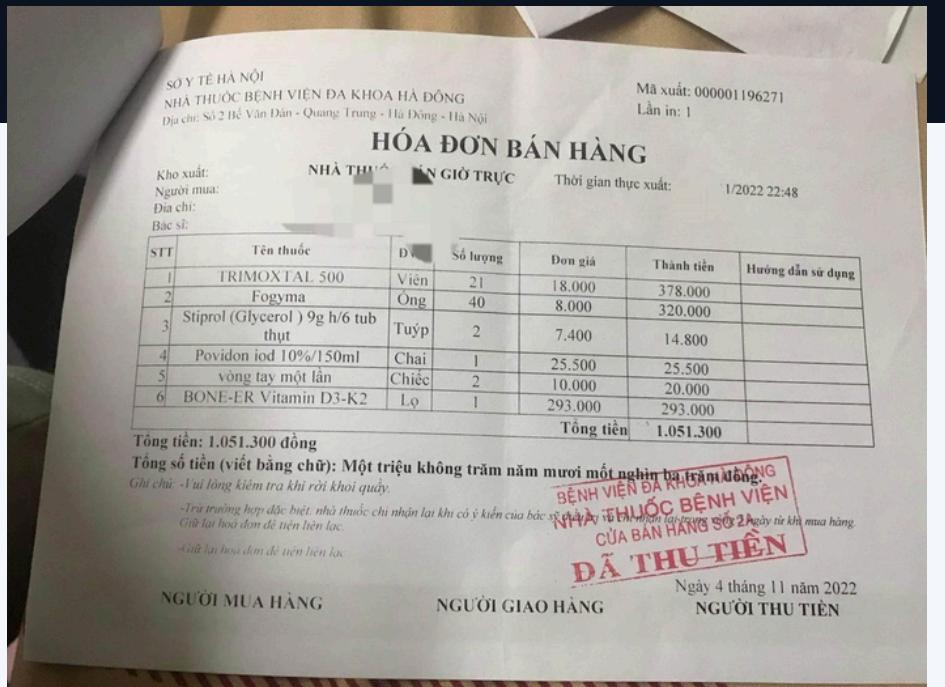
TMA Solutions

# Collect data to evaluate accuracy

## for TR model

- █ TB-BangDiem-8
- █ TB-Donthuoc-9
- █ TB-GiayDkyThue-6
- █ TB-GiayHosoXeOto-3
- █ TB-Giaykhamsuckhoe-5
- █ TB-Giaynoptien-7
- █ TB-GiayThaydoiKinhdoanh-3
- █ TB-Hopdongmuaban

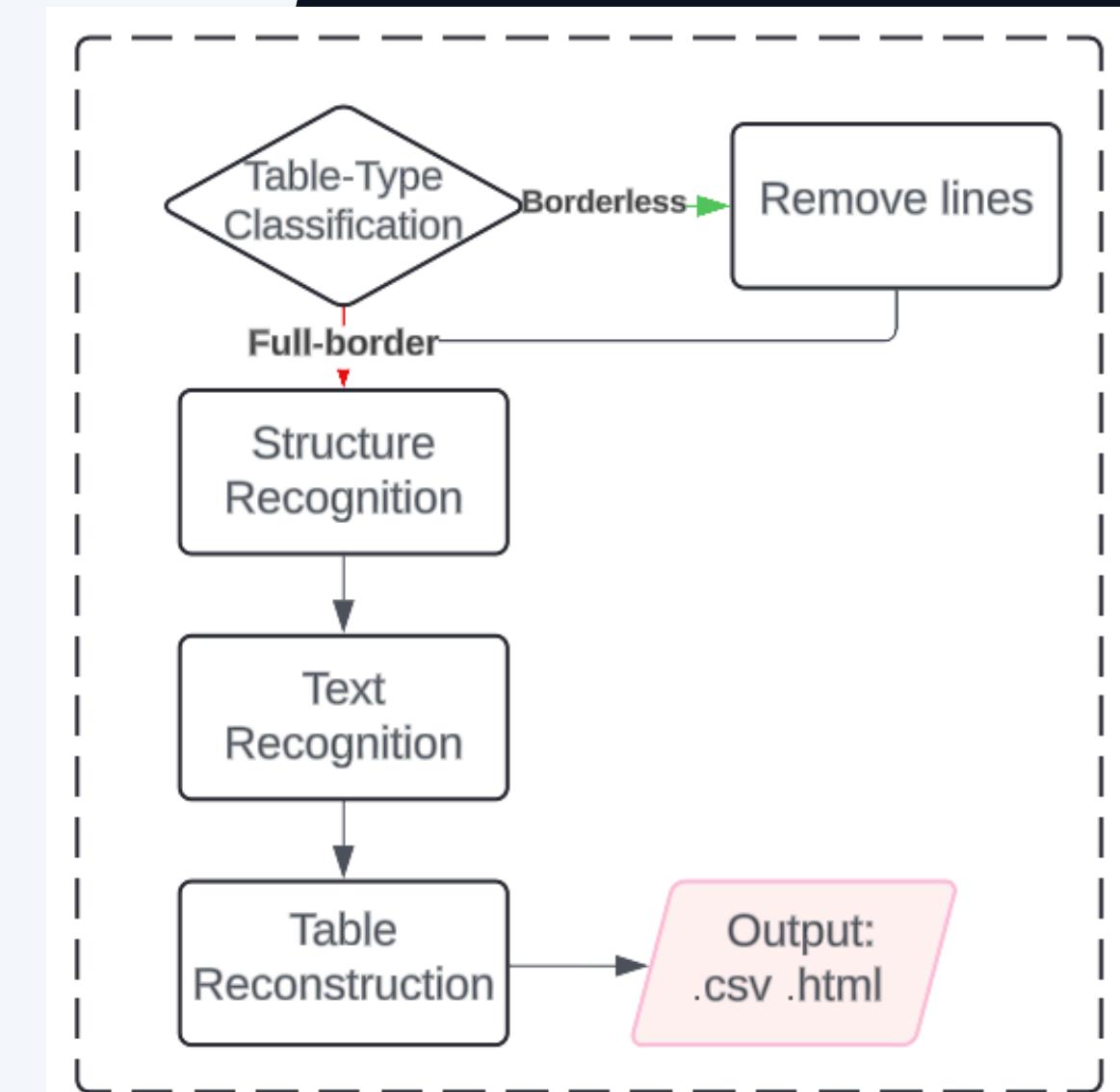
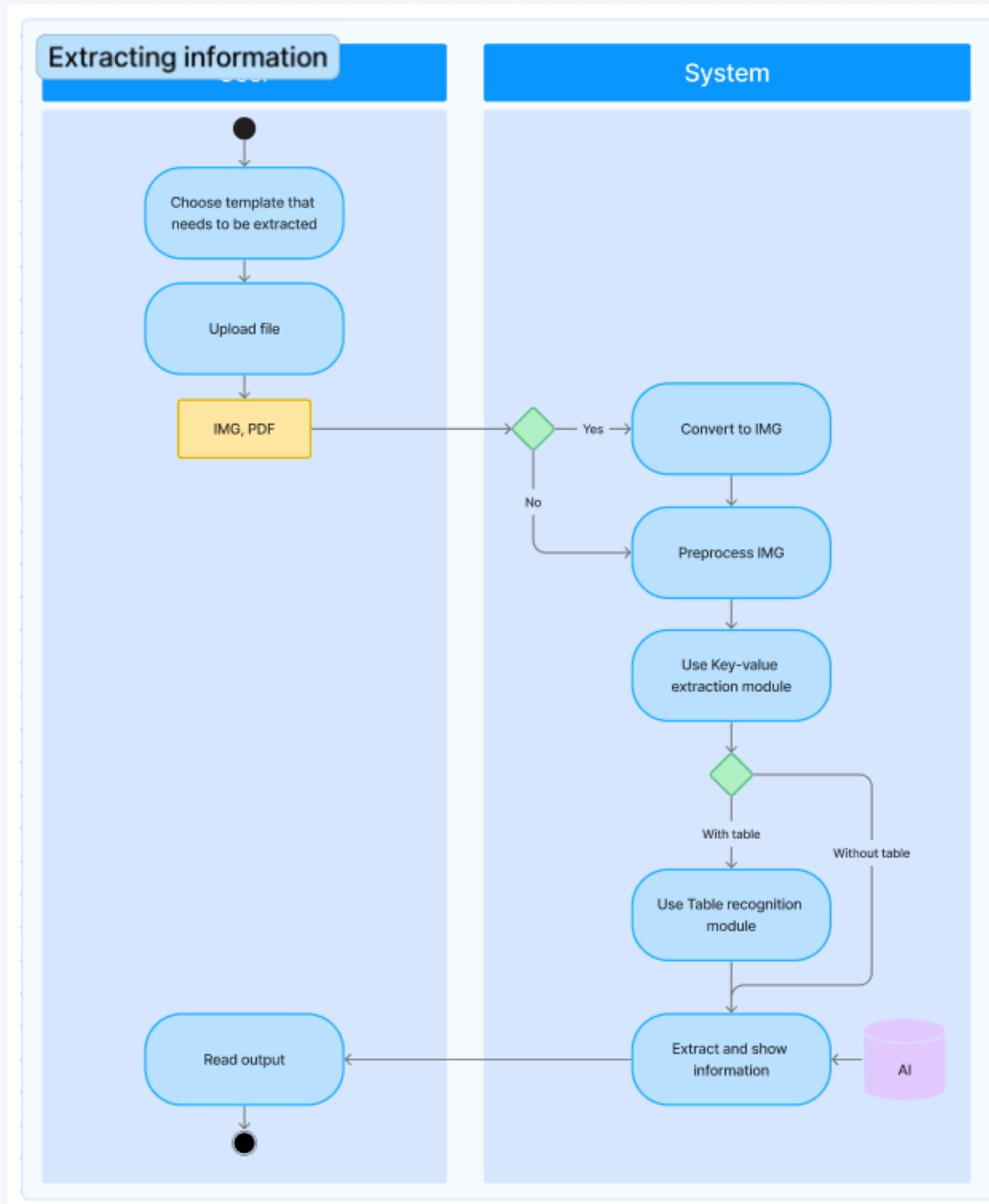
TT	Ngành đào tạo	Mã ngành	Mã tổ hợp xét tuyển	Điểm chuẩn
<b>KHÓI KINH TẾ, KỸ THUẬT, KHOA HỌC XÃ HỘI - NHÂN VĂN</b>				
1.	Kế toán	7340301	A00, A01, D01	15.0 15.5
2.	Tài chính - Ngân hàng	7340201	A00, A01, D01	15.0 15.5
3.	Quản trị kinh doanh	7340101	A00, A01, D01	15.0 15.5
4.	Quản trị văn phòng	7340406	A00, A01, D01, C00	15.0 15.5
5.	Kỹ thuật điện	7520201	A00, A01	15.0 15.5
6.	Công nghệ thông tin	7480201	A00, A01, D01	15.0 15.5
7.	Kinh tế	7310101	A00, A01, D01	15.0 15.5
8.	Ngôn ngữ Anh	7220201	D01, D09, D10, A01	15.0 15.5
9.	Marketing	7340115	A00, A01, D01	15.0 15.5
10.	Công nghệ kỹ thuật điện tử - viễn thông	7510302	A00, A01, D01	15.0 15.5



It's hard to gather paper documents that have already been taken by a camera because most of them are personal information, images collected online, and the information will be blurred, so to get enough data, the train model has to prepare its own paper so that the number of images is enough for the model to be effective.

# Draw Dynamic OCR full-flow

After the demo and feedback from the Product Manager, the workflow for the TR process will have a minor change. If the classification results in a borderless output, it will automatically fill the line and perform reconstruction. The reconstructed version will then be displayed for the user to select and make any necessary corrections, ensuring the template accurately meets the user's requirements.



# Modify & Optimize TR code

## Add Save .html .csv

```
import csv
with open('output.csv','w') as result_file:
    wr = csv.writer(result_file, dialect='excel')

    for row, row_text in data.items():
        wr.writerow(row_text)
```

Write an additional 'save' function to save it as an .html file to make it convenient for frontend developers to redraw it on the interface so that the customer can see and confirm that the template requirements have been met.

## Add Draw Table from .json

Provide a function to redraw the table, allowing the coders to visually inspect the results before and after fetching the JSON file from the model. This will help verify if the returned data matches the requirements.

```
import pandas as pd
from IPython.display import HTML
import json

with open('data.json', 'r') as f:
    data = json.load(f)
df = pd.DataFrame(data)
html = df.to_html(index=False)
display(HTML(html))
```

# Modify & Optimize TR code

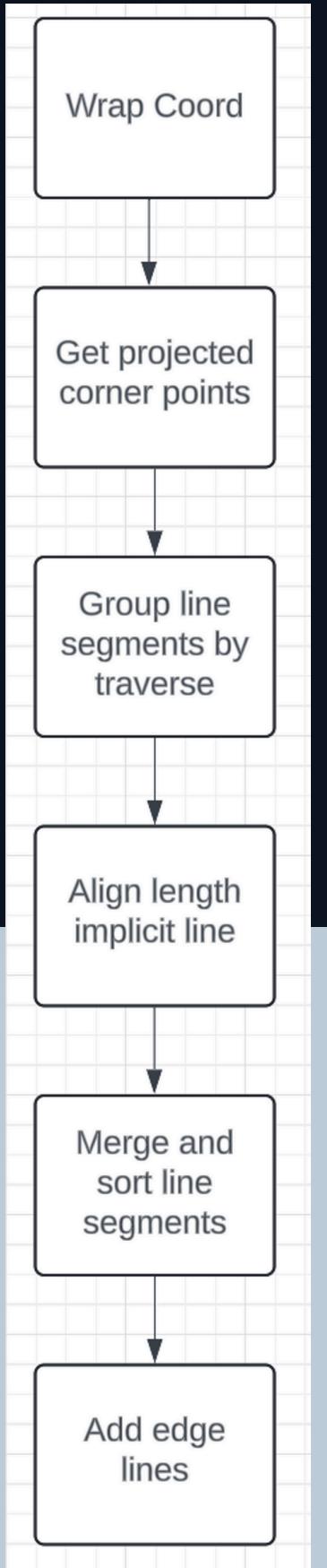
## Add Crop Cell Image & OCR Text

```

x1 = int(quad[0][0])
y1 = int(quad[0][1])
x2 = int(quad[1][0])
y2 = int(quad[2][1])
cropped_cell = image[y1:y2, x1:x2]
cv2.imwrite(f'cell_{k}.jpg', cropped_cell)
result = reader.readtext(cropped_cell)
text = "\n".join([r[1] for r in result])
text = pytesseract.image_to_string(Image.open(f'cell_{k}.jpg'),
                                    lang='vie')
text_results.append(text)
cells[k]["text"] = text_results
    
```

The model returns the coordinates of each cell and then uses those coordinates to crop the image of each cell and perform OCR text extraction using pytesseract.

However, adding this functionality makes the model run quite slowly, as it has to process multiple tasks. So I am considering whether I should include this function or just focus on reconstruction and OCR in one pass with key-value pairs.





# Collect data for training Label Medicine Box

Capture photos of the label information on a total of 17 medication bottles from different angles and backgrounds to obtain 1020 images.

Then, use Roboflow to manually detect the labels in those images, and subsequently train a YOLOv8 model to obtain the best possible model.

1020 Total Images [View All Images →](#)



Dataset Split

SET	Percentage	Count
TRAIN SET	70%	714 Images
VALID SET	20%	204 Images
TEST SET	10%	102 Images

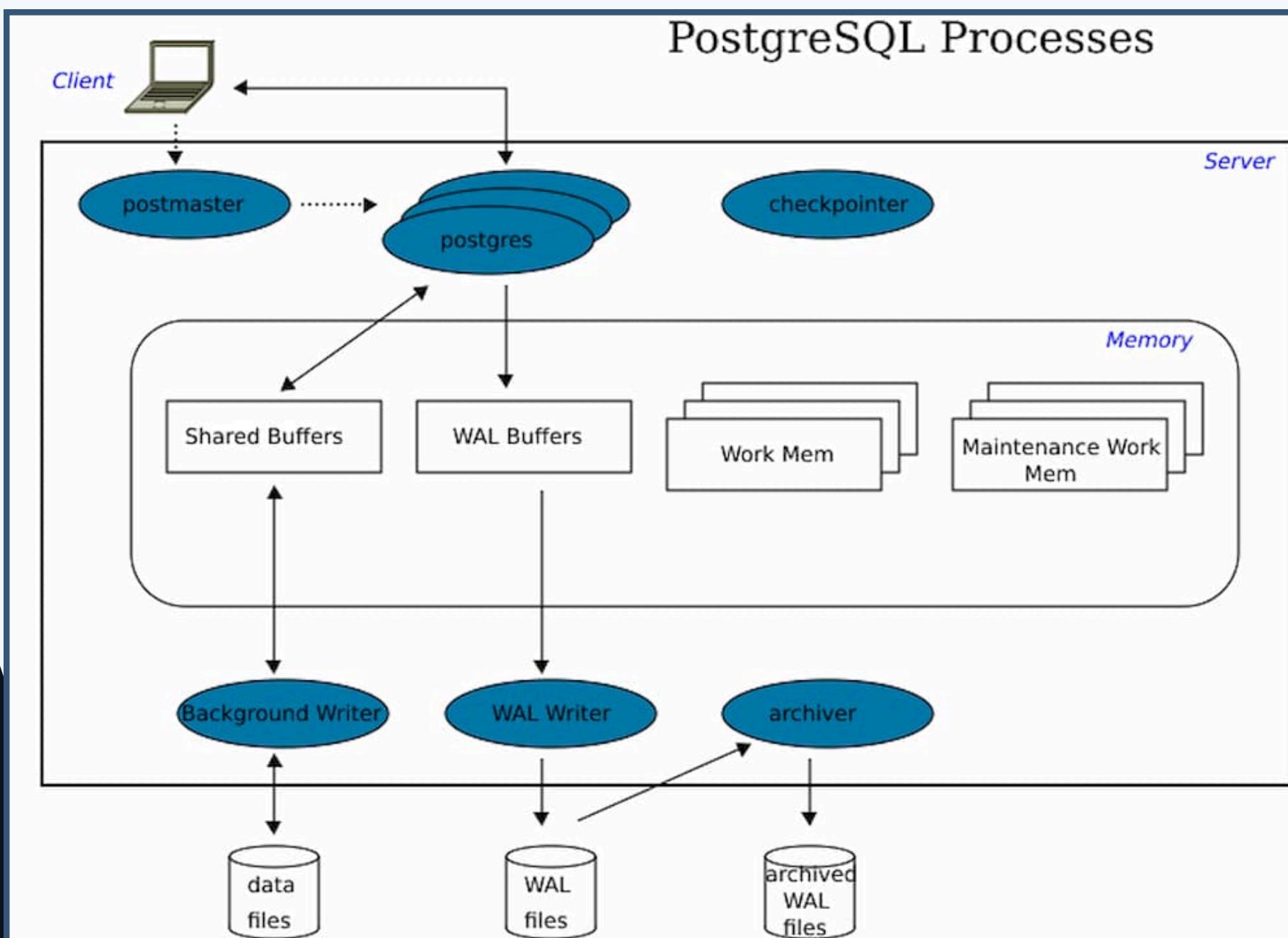
```
Ultralytics YOLOv8.2.2 🚀 Python-3.10.12 torch-2.2.1+cu121 CUDA:0 (Tesla T4, 15102MiB)
Model summary (fused): 168 layers, 11125971 parameters, 0 gradients, 28.4 GFLOPs
Downloading https://ultralytics.com/assets/Arial.ttf to '/root/.config/Ultralytics/Arial.ttf'...
100%|██████████| 755k/755k [00:00<00:00, 16.6MB/s]
val: Scanning /content/drug_bottle_detection-1/valid/labels... 204 images, 0 backgrounds, 0 corrupt: 100%|██████████| 204/204 [00:00<00:00, 1292.70it/s]
val: New cache created: /content/drug_bottle_detection-1/valid/labels.cache
```

Class	Images	Instances	Box(P)	R	mAP50	mAP50-95):	100% ██████████  13/13 [00:09<00:00, 1.38it/s]
all	204	212	0.945	0.991	0.992	0.953	

Speed: 2.6ms preprocess, 9.3ms inference, 0.0ms loss, 14.2ms postprocess per image  
Results saved to runs/detect/val2

# Research on PostgreSQL

PostgreSQL is an open-source, stable database that supports advanced SQL features. It extends SQL with scalable and reliable data handling, making it suitable for diverse applications like mobile, web, geospatial, and analytics.



## Using object-oriented features in PostgreSQL:

- Communicate with the database servers using objects in their code.
- Define complex custom data types.
- Create functions that work with their own data types.
- Define inheritance, or parent-child relationships, between tables.



# Research on

# Flask REST API



Flask API v2.0

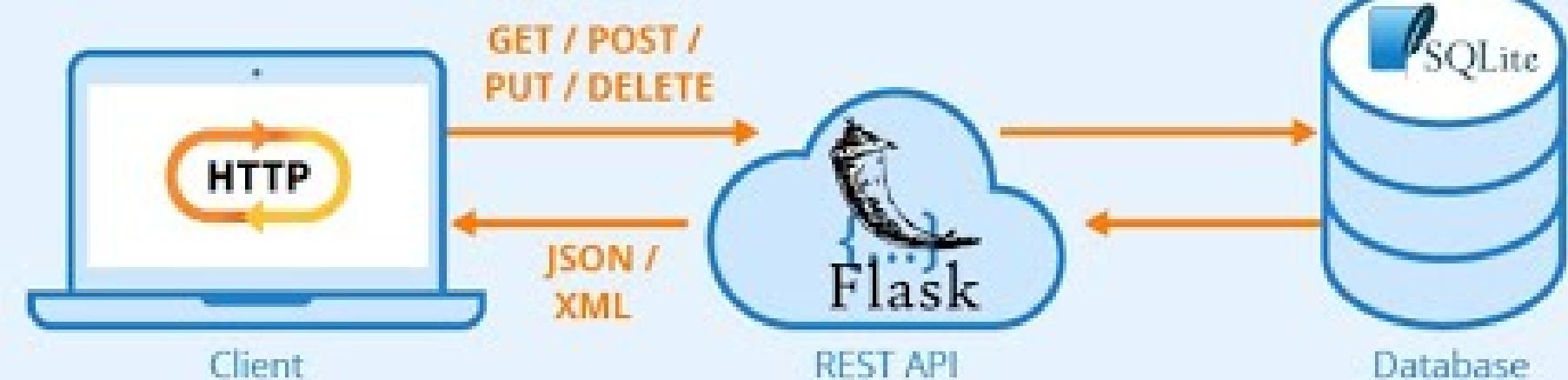
## Notes detail

Retrieve, update or delete note instances.

GET <http://127.0.0.1:5000/1/>HTTP 200 OK  
Content-Type: application/json{  
  "url": "<http://127.0.0.1:5000/1/>"  
  "text": "build the codez"  
}

Media type:

Content: `{"text": "flask-api is ace"}`



Flask is a well-known framework for web and app development, particularly for creating APIs on the Python platform.

The API is the bridge that transmits data between the client and the server, allowing the application's data to be rendered in an HTML template. The main data interaction methods on the API are GET, POST, PUT, PATCH, and DELETE, but for beginners, understanding GET and POST is typically sufficient.



# Demo TR - TRACE

localhost:7860

input\_img

Drop Image Here  
- or -  
Click to Upload

output 0

output 1

output 2

Examples

Use via API · Built with Gradio

input\_img

output 0

output 1

output 2

```
[{"id": 1, "angle": "0.3", "rect": [0: 83, 1: 481, 2: 461, 3: 221], "quad": [0: 82, 1: 479, 2: 544, 3: 483, 4: 543, 5: 784, 6: 81, 7: 781], "cells": [{"row_range": [0: 0, 1: 0], "col_range": [0: 0, 1: 0], "text": "res_input.xlsx", "confidence": 0}]}]
```

# Demo TR - Transformer

BỘ Y TẾ BỆNH VIỆN E		Mã BN:2209017539
ĐƠN THUỐC		
Họ tên: NGUYỄN VĂN THẮNG Tuổi: 39 Năm sinh: 1982 Giới tính: Nam Địa chỉ: Kim Nỗ, Đông Anh, Hà Nội, Việt Nam Chẩn đoán: K25 - Loét dạ dày / K52.3 - Viêm đại tràng không xác định		
Thuốc điều trị:		
STT	Tên thuốc - Hàm lượng	
1	Levofloxacin(Nirdicin ) - 250 mg	Uống ngày 4 viên, chia 2, sau ăn
2	Reprat 40mg(Reprat 40mg)	Uống ngày 2 viên, chia 2 lần - 6h - 21h
3	L-glutathion 1000mg , Collagen natri 400mg...(Eu- Thion )	Uống ngày 2 viên, chia 2, sau ăn
4	Dầu gan cá mập 250mg(Ecomer )	Uống ngày 2 viên, chia 2, sau ăn
5	L-Carnitine 200mg(B12 Energy )	Ngày uống 1 lọ, sau ăn
6	(Intesta) - 500 mg	Uống ngày 3 viên, chia 3 lần - sau ăn
Cộng khoán: 6		

```

Max number of columns: 2
['STT', 'Tên (uuC']
['Levofloxacin(Nirdicin ) - 250 mg [Uống ngày 4 viên, chia 2. sau ăn', '']
['Reprat 40mg(Reprat 40mg) [Uống ngày 2 viên chia 2 lần 6h - 21h', '']
[['L-glutathion 1000mg , Collagen natri 400mg- (Eu- Thion Uống ngày 2 viên, chia 2 sau &n', '']
[['Dầu gan cá mp 250mg(Ecomer luóng ngày 2 viên; chia 2, sau &n', '']
['L-Camitinc 200mg(B12 Energy [Ngày uong 1o, sau ún', '']
['(Intesta ) 500 Uống ngày 3 viên; chia 3 lan sa...', '']

```

BỘ Y TẾ BỆNH VIỆN E		Mã BN:2209017539
ĐƠN THUỐC		
Họ tên: NGUYỄN VĂN THẮNG Tuổi: 39 Năm sinh: 1982 Giới tính: Nam Địa chỉ: Kim Nỗ, Đông Anh, Hà Nội, Việt Nam Chẩn đoán: K25 - Loét dạ dày / K52.3 - Viêm đại tràng không xác định		
Thuốc điều trị:		
STT	Tên thuốc - Hàm lượng	
1	Levofloxacin(Nirdicin ) - 250 mg	Uống ngày 4 viên, chia 2, sau ăn
2	Reprat 40mg(Reprat 40mg)	Uống ngày 2 viên, chia 2 lần - 6h - 21h
3	L-glutathion 1000mg , Collagen natri 400mg...(Eu- Thion )	Uống ngày 2 viên, chia 2, sau ăn
4	Dầu gan cá mập 250mg(Ecomer )	Uống ngày 2 viên, chia 2, sau ăn
5	L-Carnitine 200mg(B12 Energy )	Ngày uống 1 lọ, sau ăn
6	(Intesta) - 500 mg	Uống ngày 3 viên, chia 3 lần - sau ăn
Cộng khoán: 6		

Table Table (rotated)

STT	Tên thuốc - Hàm lượng
1	Levofloxacin(Nirdicin ) - 250 mg
2	Reprat 40mg(Reprat 40mg)
3	L-glutathion 1000mg , Collagen natri 400mg...(Eu- Thion )
4	Dầu gan cá mập 250mg(Ecomer )
5	L-Carnitine 200mg(B12 Energy )
6	(Intesta) - 500 mg

STT	Tên thuốc - Hàm lượng
0	Levofloxacin(Nirdicin ) - 250 mg [Uống ngày 4 ...
1	Reprat 40mg(Reprat 40mg) [Uống ngày 2 viên chia ...
2	[L-glutathion 1000mg , Collagen natri 400mg- (...
3	[Dầu gan cá mp 250mg(Ecomer luóng ngày 2 viên; ...
4	L-Camitinc 200mg(B12 Energy [Ngày uong 1o, sau ún
5	(Intesta ) 500 Uống ngày 3 viên; chia 3 lan sa...



# Review Demo TRACE

## Advantages:

- Fast, Quite high accuracy, Can detect span-cells
- Separate modules for two classes - fullborder/ borderless

## Disadvantages:

- Even for borderless tables, it reconstructs them as fullborder
- Tends to merge two close tables into one
- Requires CUDA usage

## Comparison to other modules:

- Faster and more accurate than modules like PaddleOCR-Transformer
- GPU usage is mandatory

STT	Tên thuốc - Hàm lượng
1	Levofloxacin( Nirdicin ) - 250 mg Uống ngày 4 viên, chia 2, sau ăn
2	Reprat 40mg(Reprat 40mg) Uống ngày 2 viên, chia 2 lần - 6h - 21h
3	L-glutathion 1000mg , Collagen natri 400mg...(Eu- Thion ) Uống ngày 2 viên, chia 2, sau ăn
4	Dầu gan cá mập 250mg(Ecomer ) Uống ngày 2 viên, chia 2, sau ăn
5	L-Carnitine 200mg(B12 Energy ) Ngày uống 1 lọ, sau ăn
6	Intesta) - 500 mg Jóng ngày 3 viên, chia 3 lần - sau ăn

Cộng khoản: 6

# Lessons learnt

I have learned the following skills:

- How to create a plan for the week and for each sprint.
- How to fill out a daily report to ensure the progress of work.
- How to share updates during each sprint review.
- How to work effectively in a team and ensure the efficiency and results of my own tasks.
- The process of researching a new topic and implementing it according to the project requirements.
- How to use GitLab to integrate the team's code.
- How to support each other to ensure the project's progress.
- Gained a lot of new knowledge from the mentor PM and colleagues.
- Received training on various topics such as the product development process and teamwork.

These skills have helped me develop a structured approach to project management, collaboration, and continuous learning to meet the demands of the project.

# Best practices

I have improved the following skills:

- Planning and executing tasks more effectively within the project.
- Enhancing my technical expertise and domain knowledge.
- Improving my ability to share problems and project updates with others.
- Becoming more extroverted, friendly, and caring towards my colleagues.
- Enhancing my English proficiency by working in an English-speaking environment.
- Consistently meeting deadlines and ensuring timely completion of assigned tasks.
- Being more proactive in communicating my issues to supervisors and co-workers to receive timely assistance.

The improvements I've made have made me a more effective and valuable team member, greatly contributing to the project's success.



# Thanks

---

Thank you to the company for providing me the opportunity to intern in the field I'm passionate about pursuing. I'm grateful to the interns in the department who have always supported us. Their help enabled me to complete the internship in the best possible way - this has been a very meaningful period for me.

I would like to express my sincere gratitude to Nguyen Minh Nghia and Tran Trung Kien for guiding and supporting me, for giving me the chance to gain experience, develop myself, and enhance my domain knowledge. I'm also thankful to my colleagues who have been by my side, supporting each other in our work and personal lives.

