

Thai  
Optical Character  
Recognition

# Open-Source Initiative Project

Last Updated: 11 Jun 2023

# เป้าหมายหลัก

1. ต้องการให้เห็นถึงความสำคัญของ open-source software ในประเทศไทย
2. ต้องการขับเคลื่อนให้ทุกปัจเจกในประเทศไทยทำ digitization และสามารถทำ data analytics ขึ้นต่ำได้โดยไม่มีค่าใช้จ่าย
3. ความสามารถในการเข้าถึงความรู้ของประชาชนไทยต้องเท่าเทียมเสมอกัน
4. ทำให้ภาครัฐสามารถทำ data transformation โดยไม่มีความจำเป็นต้องจ้างจากเอกชน ซึ่งสามารถป้องกันเรื่องข้อมูลความมั่นคงและการขัดกันทางผลประโยชน์

ความฝัน  
อันสูงสุด

## Data-driven nation

Digitization for everyone

Data analytics for everyone

Open data-information-knowledge-wisdom  
for everyone



Thai Optical Character Recognition

Open-Source  
Initiative Project

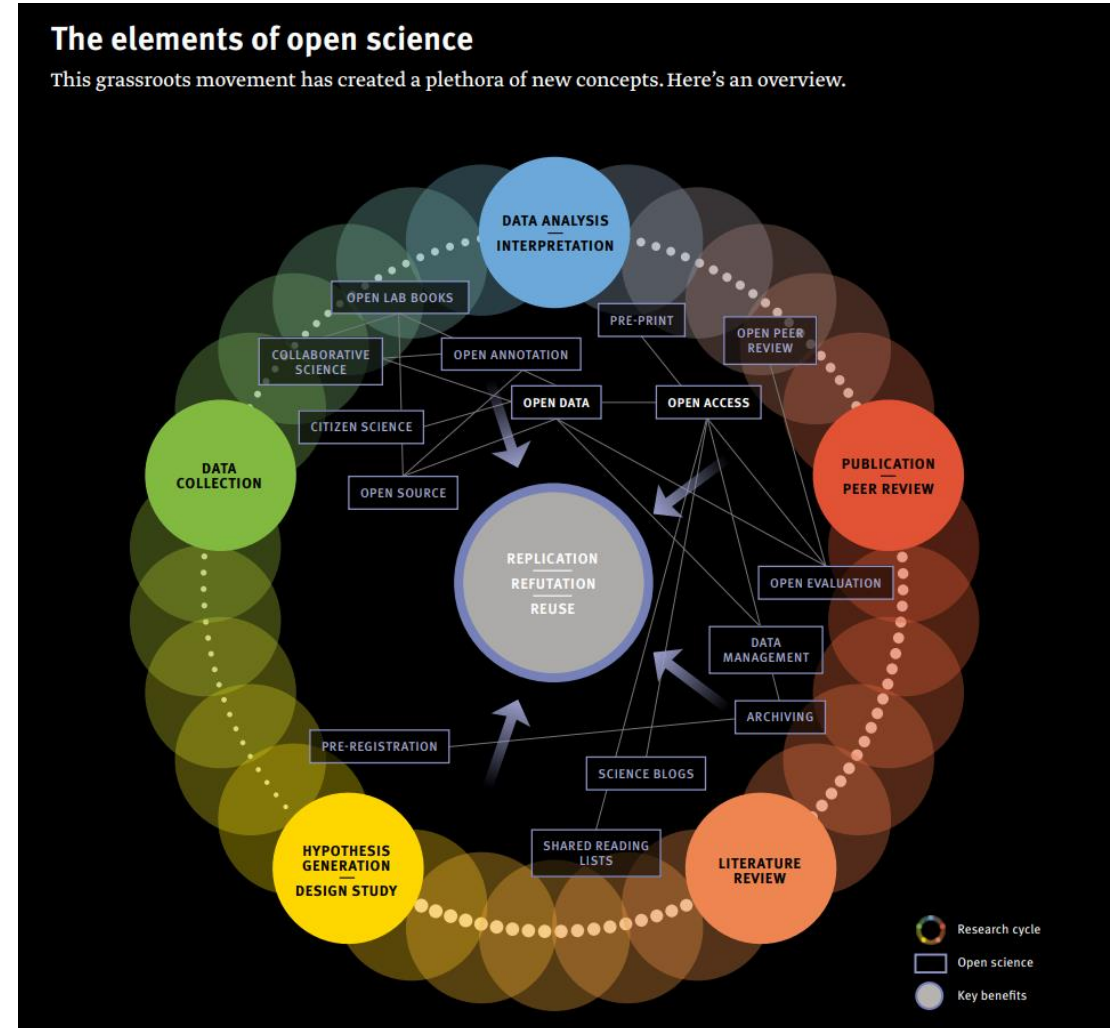
# ความพยายามในเป้าหมายนั้น

## My approach

1. Thai OCR engine open-source initiative
2. Thai open repository for advanced studies (more detailed than Wikipedia Thai; beyond the language barrier)

## Governmental approach

1. Open National Data
2. Big Data and AI-related governmental organization
3. Smart city



Thai Optical Character Recognition

Open-Source  
Initiative Project

# Market Analysis

## Current Progress

## Future Direction

# Business Model

## Key partners

Existing open-source models  
(Tesseract-OCR, PyThaiNLP, OpenCV)

## Key activities

- Friendly OCR tools
- Knowledge graph construction over documents

## Customer segments

(Prioritizing from most to least)

- 
- Governmental organizations

- 
- SMEs

- 
- Individuals who'd like to have personal archives

## Cost structure

- Man-day costs
- Data analytics tools and pricing for renting cloud services (maybe)
- Specialist consults/lectures

## Value propositions

General

- Easier digitization
- Easier knowledge discovery process
- **"Always" open-source tools for everyone**

- 
- Reducing processing time for data archiving
  - Open-data initiation

- 
- Importing data for data analytics process
  - Precision accounting
  - Data-driven organization in no-price

- 
- Easy-to-archive

## Revenue Streams

This project is for charity.  
NO ANY DIRECT INCOME.

Indirect income sources:

- Open Science
- Open National Data Initiation

## Customer Relationships

- Support Community in GitHub
- Bug report and feature proposal

## Channels

- Developer society
- GitHub
- Governmental/academic organization (for free tools)
- Release notes

Thai Optical Character Recognition

Open-Source  
Initiative Project

## การการันตีกลุ่มลูกค้า

### องค์กรของรัฐ

- เอกสารมีความหลากหลายและเก๋าจำนวนมาก
- ความจำเป็นในการวิเคราะห์ข้อมูลและนโยบายของรัฐที่พยายามผลักดันให้เกิด data-driven organization
- ประวัติขององค์กรและการตามเอกสาร

### ตัวอย่างองค์กร

- สนง.คกก. ฤๅษฎีภา
- หจช. และ หสช.
- หน่วยงานภาครัฐอื่น ๆ ที่ต้องการทำ record management systems
- SMEs
- บุคคลผู้สนใจงานจดหมายเหตุและงานทางโบราณคดี

## ความเข้าใจที่ต้องมี

### องค์กรของรัฐ

- หน่วยงานต้องการทราบว่าตัวเองถือข้อมูลอะไรอยู่
- หน่วยงานต้องการทำ data analytics กับเอกสารที่มี
- หน่วยงานมีภารกิจที่ต้องปล่อยข้อมูลลง GD Catalogue, GDX และระบบสารบรรณอิเล็กทรอนิกส์อื่น ๆ

### คุณภาพของข้อมูล

- ต้องสามารถวิเคราะห์ได้โดยง่ายและไม่คลุมเครือ (5-star open data <https://5stardata.info/en/>)
- ต้อง digitize ทุกรายละเอียดของเอกสาร การมี metadata ที่ดีเป็นสิ่งสำคัญ
- Versioning เป็นสิ่งจำเป็นในการเก็บเอกสาร

Market Analysis

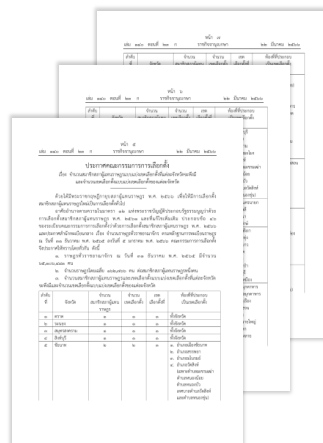
**Current Progress**

Future Direction

# ภาพรวมของ OCR

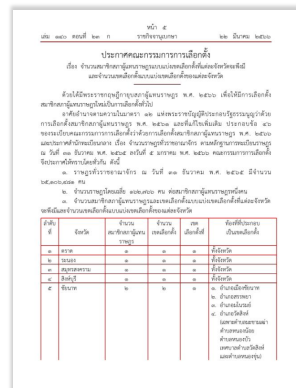


ไฟล์ PDF



ไฟล์ภาพ lossless (PNG)

ตรวจจับเส้นตาราง  
(Straight Line Detection)



ตรวจจับกลุ่มตัวอักษร  
(Region of interest detection)



แปลงตารางเป็น CSV



สร้างโครงสร้างเป็น  
Markdown

รู้จำตัวอักษร  
(Text recognition)



Tesseract OCR

Customized  
Tesseract  
Engine

๓. จำแนกกันทวิชัย

๓. จำแนกกันทวิชัย

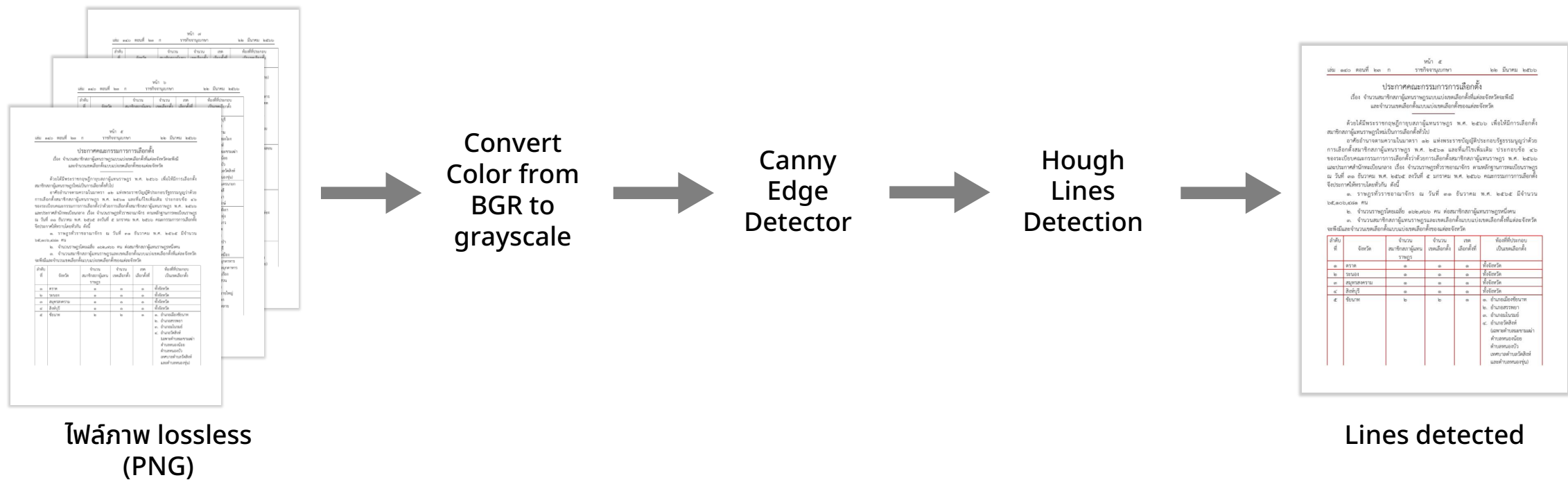
Text region

Thai Optical Character Recognition

Open-Source  
Initiative Project

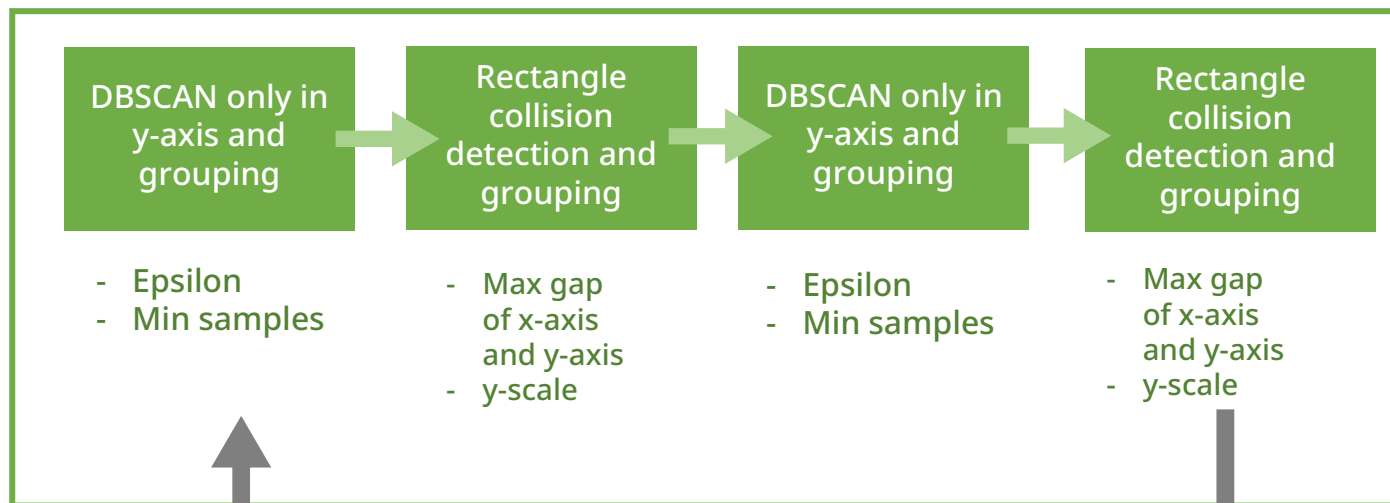


# Straight line detection



# Region of interest detection

Hyperparameter auto-tuning:  
using Grid search for optimization



Objective  $\min_{\theta} L(\theta)$

$$L(\theta) = \sum_i \text{Area}(\text{Rect}(\theta, i))$$

สิ่งที่ต้องการใน optimization คือ  
ทำให้พื้นที่กล่องที่ใช้คลุมตัวอักษร  
จะต่อน้อยที่สุด เพื่อป้องกันการคลุม  
กล่องในหลายบรรทัด

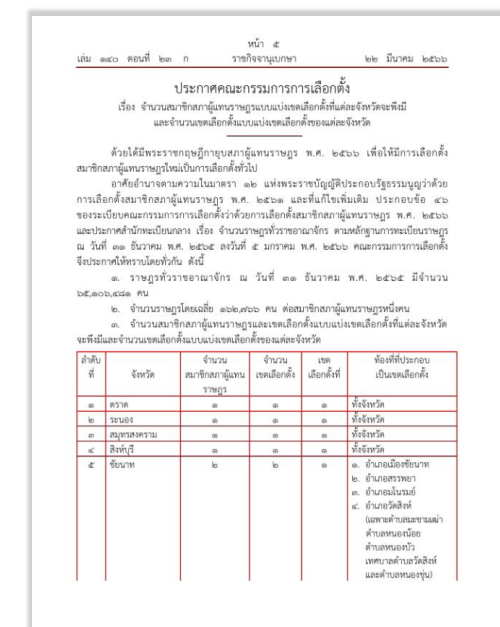
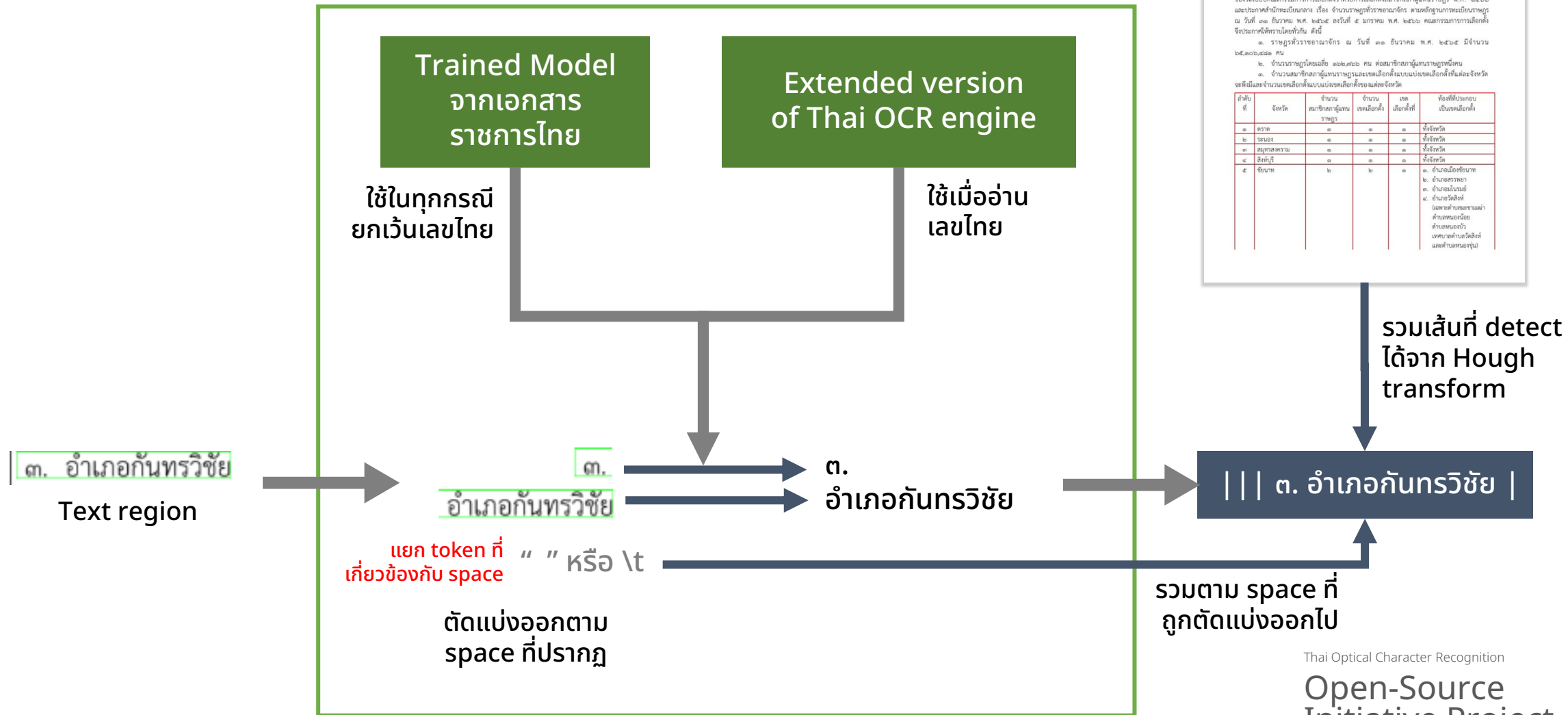
ตรวจสอบกลุ่ม  
ตัวอักษร  
(Region of  
interest  
detection)

x1	y1	x2	y2	text
100	100	200	150	๑. ๑๒๓๔๕๖๗๘๙
100	150	200	200	๒. ๑๒๓๔๕๖๗๘๙
100	200	200	250	๓. ๑๒๓๔๕๖๗๘๙
100	250	200	300	๔. ๑๒๓๔๕๖๗๘๙
100	300	200	350	๕. ๑๒๓๔๕๖๗๘๙
100	350	200	400	๖. ๑๒๓๔๕๖๗๘๙
100	400	200	450	๗. ๑๒๓๔๕๖๗๘๙
100	450	200	500	๘. ๑๒๓๔๕๖๗๘๙
100	500	200	550	๙. ๑๒๓๔๕๖๗๘๙
100	550	200	600	๑๐. ๑๒๓๔๕๖๗๘๙

Thai Optical Character Recognition

Open-Source  
Initiative Project

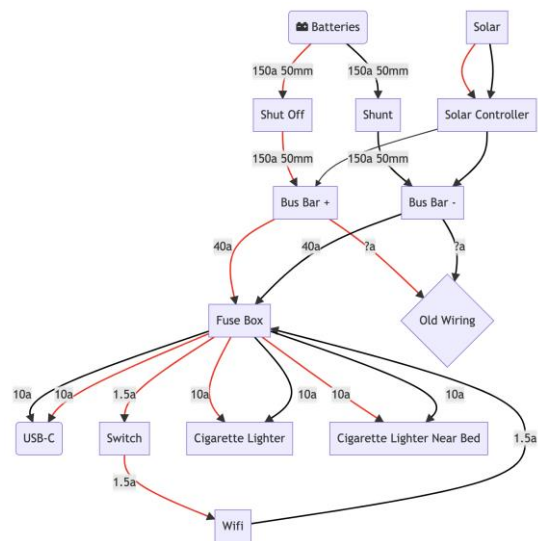
# Text recognition



Market Analysis  
Current Progress  
**Future Direction**

# Pain points

Diagram export from OCR engine



We will format it as HTML/Markdown format.

More fine-tuning of OCR model

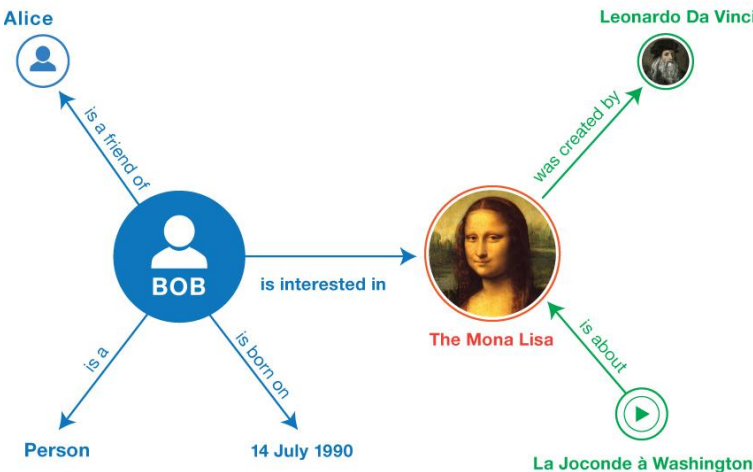


Tesseract OCR

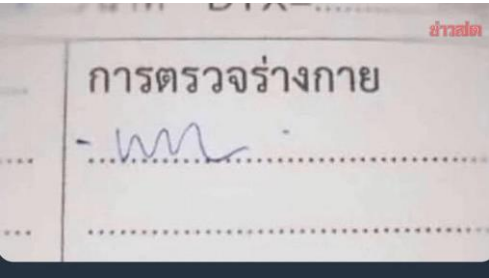
Variety of textual communication



Knowledge graph construction from metadata which is defined in the document



- Quick-win project
- Metadata generation



Handwriting/Typing classification

Thai Optical Character Recognition

Open-Source Initiative Project

# The quick-win project

ชุดที่ ๑ (สีขาว) ส่งให้คณะกรรมการการเลือกตั้ง  
ชุดที่ ๒ (สีเขียว) ปิดประทับ  
ชุดที่ ๓ (สีชมพู) ใส่ในถุงวัสดุใสขึ้นนอก

ส.ส. ๕/๑๘

รายงานผลการนับคะแนนสมาชิกสภาผู้แทนราษฎรแบบแบ่งเขตเลือกตั้ง

ตามที่ได้มีพระราชกฤษฎีกาให้มีการเลือกตั้งสมาชิกสภาผู้แทนราษฎร และคณะกรรมการการเลือกตั้ง ได้กำหนดให้วันที่ ๑๔ เดือน พฤษภาคม พ.ศ. ๒๕๖๖ เป็นวันเลือกตั้ง นั้น

บัดนี้ คณะกรรมการการประจำหน่วยเลือกตั้งได้ดำเนินการนับคะแนนสมาชิกสภาผู้แทนราษฎร แบบแบ่งเขตเลือกตั้งของหน่วยเลือกตั้งที่ ๕๑ หมู่ที่ ๑ ตำบล/แขวง/เทศบาล (สีน้ำเงิน)

อำเภอ/เขต (สีน้ำเงิน) เขตเลือกตั้งที่ ๕ จังหวัด กรุงเทพมหานคร เสร็จสิ้นเป็นที่เรียบร้อยแล้ว

ดังนั้น จึงขอรายงานผลการนับคะแนนของหน่วยเลือกตั้งดังกล่าว ดังนี้

๑. จำนวนผู้มีสิทธิเลือกตั้ง

๑.๑ จำนวนผู้มีสิทธิเลือกตั้งตามบัญชีรายชื่อผู้มีสิทธิเลือกตั้ง 4๐๐ คน (สี่ร้อยคน)

๑.๒ จำนวนผู้มีสิทธิเลือกตั้งที่มาแสดงตน 292 คน (สองร้อยเก้าสิบสอง) (เฉพาะวันเลือกตั้ง)

21 พ.ค. 2565

cr.กกต.

เข้าไปตรวจสอบกัน  
กกต.เปิดคะแนนเลือกตั้งรายหน่วย  
ให้ดูในเว็บไซต์แล้ว

## แบบ ส.ส. 5/18

- มีจำนวน ~95,000 หน่วยเลือกตั้ง
- ไฟล์อยู่ในรูปแบบ scanned PDF
- ขนาดไฟล์รวม ~320 GB

## Tasks

- Metadata generation
- Directory tree rearrangement
- Digitization process

## Opportunities

- New datasets for Thai OCR
- Large-scale digitization prototyping

Thai Optical Character Recognition

Open-Source  
Initiative Project

# Assignments

## The quick-win project

Project manager: บิว

Dev: (ตามความเหมาะสม)

## Tasks ที่ต้องทำก่อนเริ่มงาน

1. Epic planning
2. Sprint backlogging
3. Man-day estimation
4. The protocol for work management
5. Repository management
6. Wiki management
7. Sprint review scheduling