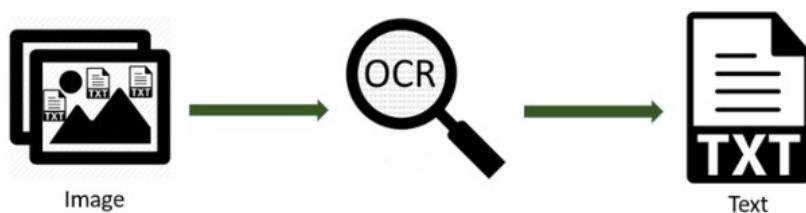


Optical Character Recognition: **Converting** text in an image into a machine-readable format



-  OCR widely studied for high-resource languages like English and Chinese
-  **Limited understanding** on effectiveness and limitations of OCR on **Southeast Asian (SEA) languages**
 - Low-resource: Less training data
 - SEA population: 700 million people

Studying OCR for SEA languages is **valuable**.

- Enhanced **digital accessibility** for low-resource language users
- **Resource creation** for downstream Natural Language Processing tasks
 - E.g., machine translation, named-entity recognition

Benchmarking and Improving OCR Systems for Southeast Asian Languages

Project ID: H0792230

Supervisor: A/P Min-Yen Kan

Advisor: Tongyao Zhu

Student Name: Qiu Jiasheng, Jason

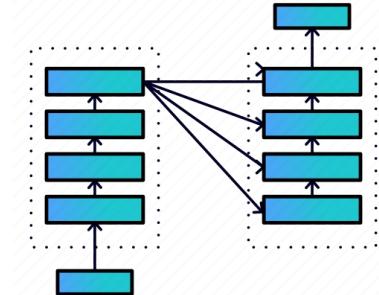
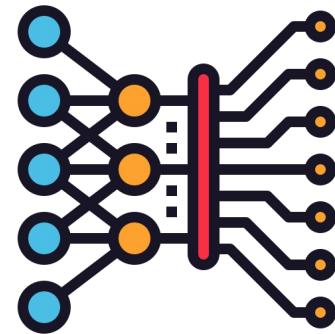
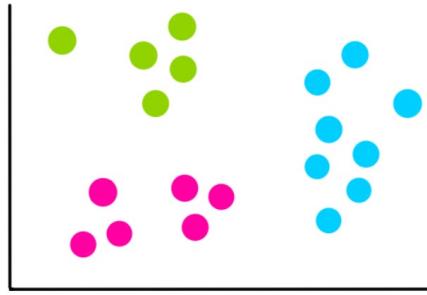
RQ1: How do popular OCR tools **perform** on SEA scripts?

RQ2: What **script-related challenges** affect OCR accuracy on SEA languages?

RQ3: How can fine-tuning **enhance** OCR accuracy on SEA languages?

- 01 Motivation
- 02 Related Work
- 03 Experiment 1: Benchmarking on
Real-world Wikipedia **Screenshots**
- 04 Experiment 2: Benchmarking on
Synthetic Data
- 05 Experiment 3: **Fine-tuning**

OCR model architectures are **evolving**.



Traditional machine learning (KNN, SVM) on manual features¹

Deep learning (CNN, RNN)¹

Transformer-based architecture²

Our research: Compares performance of **different OCR models** on SEA languages

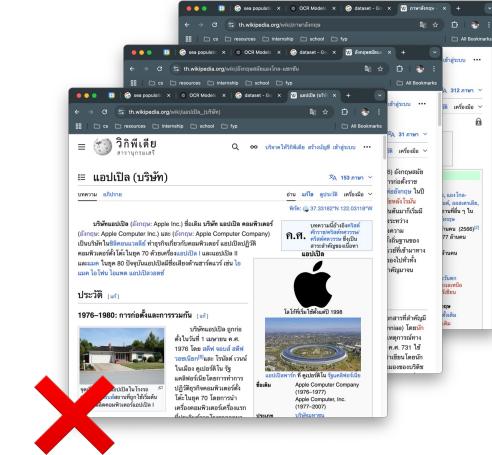
¹ A survey of deep learning approaches for OCR and document understanding (Subramani et al., 2020)

² Fine-tuning vision encoder-decoder transformers for handwriting text recognition on historical documents (Parres & Paredes, 2023)

OCR on SEA languages is **challenging**.

ทุกคนมีลิทธิ์ที่จะออกจากประเทศไทย ๆ ไป
รวมทั้งประเทศของตนเองด้วย และที่จะ^{จะ}
กลับยังประเทศไทย

SEA scripts are more
complex



Lack of quality datasets

Our research: Develops a reusable **pipeline** for **collecting**
real-world data and **generating** synthetic data

Synthetic data with noise helps with data scarcity.

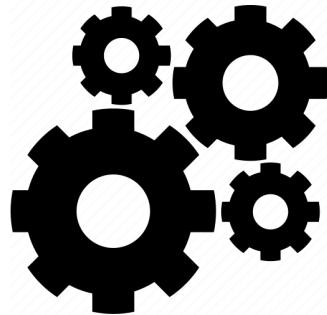


Easy to generate large-scale, quality data for low-resource languages

Adding **noise** better simulates real-world conditions

Our research: Adopts **similar** techniques by generating and benchmarking on synthetic data with noise

Fine-tuning can adapt OCR models to new languages.



Fine-tuning: **Further training** pre-trained models on task-specific dataset

Retains prior knowledge +
Adapts to new dataset



Limited data can achieve **good** results

Our research: **Fine-tunes** existing model on SEA languages

01

Motivation

02

Related Work

03

Experiment 1: Benchmarking on
Real-world Wikipedia **Screenshots**

04

Experiment 2: Benchmarking on
Synthetic Data

05

Experiment 3: **Fine-tuning**

Experiment **setup**

Experiment 1: Benchmarking on
Real-world Wikipedia **Screenshots**

Experiment 2: Benchmarking on
Synthetic Data

Experiment 3: **Fine-tuning**



1. Methodology
2. Results
3. Discussion

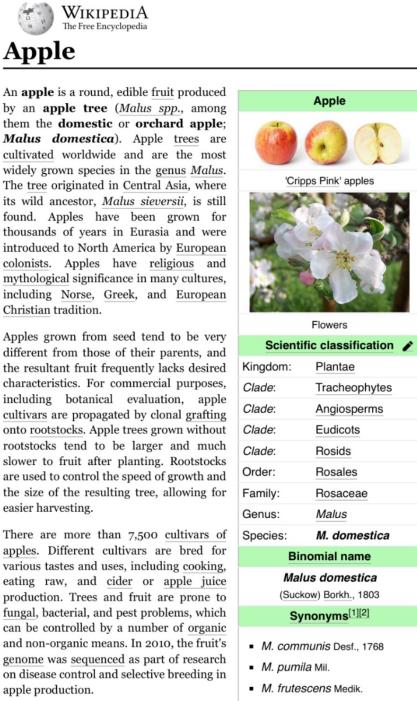
4 languages: English, Indonesian, Vietnamese, Thai

	Speaker population	Script Type	Example
English	1.5 billion	Latin	Good morning
Indonesian	252 million	Latin	Selamat pagi
Vietnamese	97 million	Latin with diacritics	Chào buổi sáng
Thai	71 million	Brahmic	สวัสดีตอนเช้า

Why?

1. Capture diverse range of **script types**
2. Wide usage makes it feasible to obtain **data**

Data source: [Wikipedia](#)



Why?

1. Can be **easily** converted into images via screenshots + **Accessible** API for article text
 2. Large **availability** of data in SEA languages

OCR systems: **EasyOCR**, **Tesseract**, **GOT**

Criteria:

1. **Different** underlying architectures
 2. **Open-source** with wide support for SEA languages
-

EasyOCR

- Text detection: **CRAFT** algorithm (CNN)
- Text recognition: **Convolutional Recurrent Neural Network**
- Supports **83** languages

Tesseract

- One of the most **well-known** open-source OCR engines
- Long Short-Term Memory (**LSTM**) RNN
- Supports **116** languages

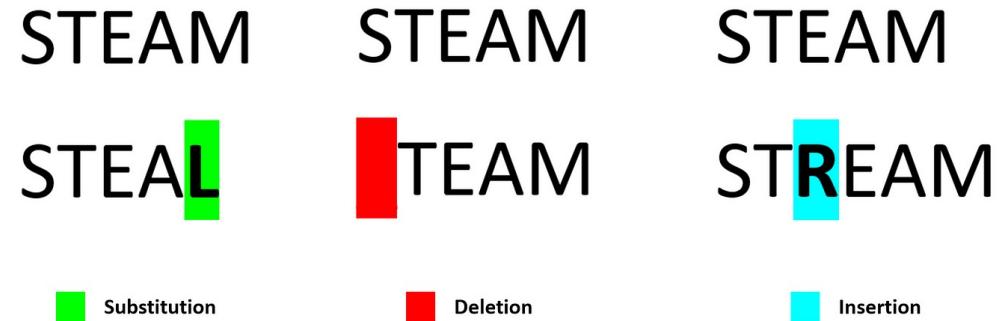
GOT
General OCR
Theory

- Designed to recognize **beyond** traditional text (e.g., sheet music)
- **Vision Encoder Decoder** architecture
- **580 million** parameters
- Supports only English and Simplified Chinese

Evaluation metrics: **CER** and **WER**

Character Error Rate (CER):

$$CER = \frac{S + D + I}{N}$$



- N: Number of characters in ground truth
- Lower value → **Better**
- **Word Error Rate (WER)**: Word-equivalent of CER
 - Commonly used when preserving entire word is important

01

Motivation

02

Related Work

03

Experiment 1: Benchmarking on
Real-world Wikipedia **Screenshots**

04

Experiment 2: Benchmarking on
Synthetic Data

05

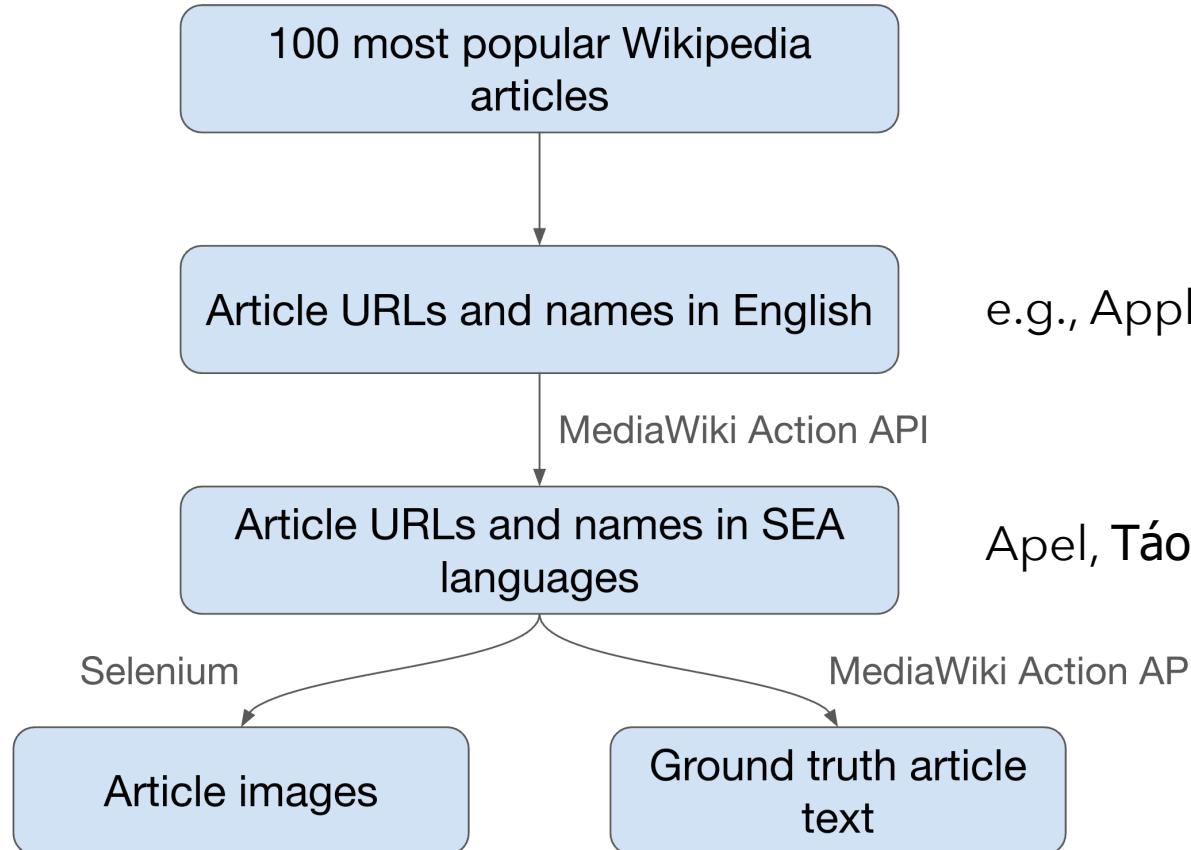
Experiment 3: Fine-tuning

Goal: Explore performance of OCR tools on SEA scripts (RQ1)

Steps:

1. **Collect** article screenshots and texts from Wikipedia
2. **Evaluate** OCR tools

1. **Collect** article screenshots and texts from Wikipedia



e.g., Apple

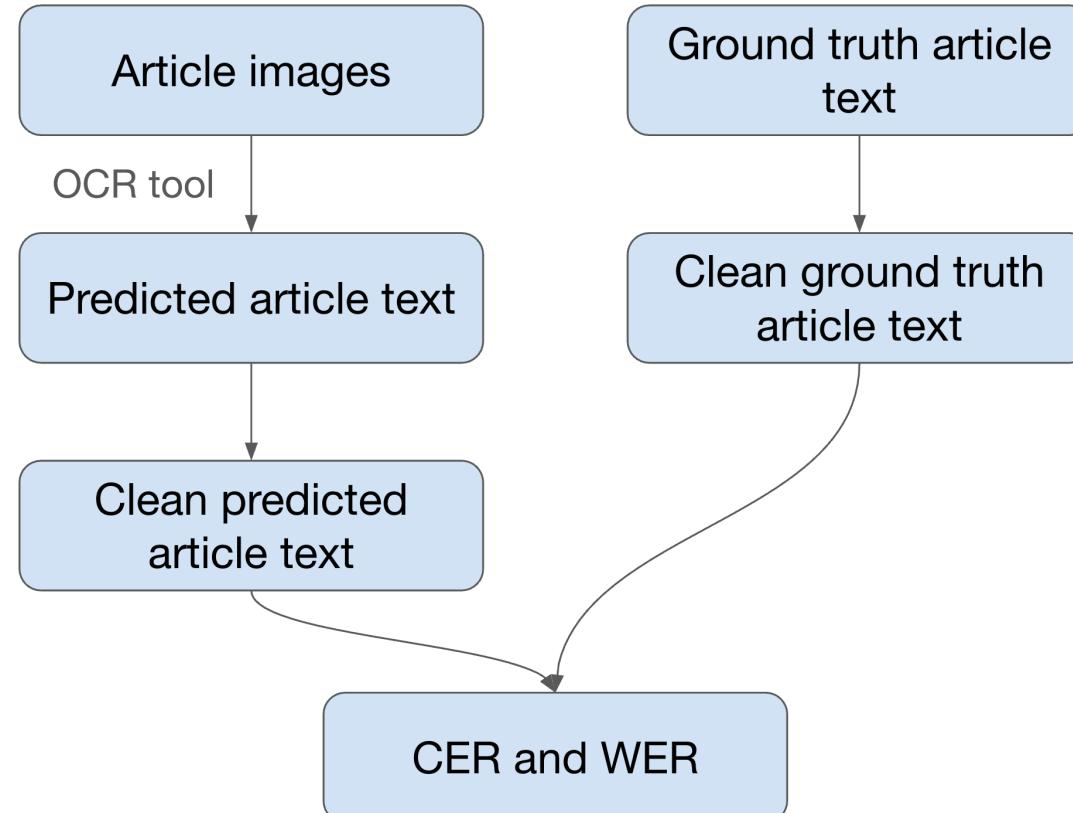
Apel, Táo tây, แอป



Result: **7,976** Wikipedia screenshots +
Ground truth text



Predicted text



CER: 0.06
WER: 0.10

ได้ มาจากต้นแอปเปิล (Malus spp.)
รวมถึง แอปเปิลบ้าน และแอปเปิลสวน
ผลไม้ (Malus domestica) ต้นแอปเปิล
ได้รับการเพาะปลูกทั่วโลก และเป็นสาย^{พันธุ์}ที่ปลูกมากที่สุดในสกุล Malus มีถิ่น^{กำเนิด}ในเอเชียกลาง โดยมี Malus sieversii เป็นบรรพบุรุษป่า ซึ่งยังคงพบ^{อยู่}ในภูมิภาคเอเชียกลาง เช่น จีน มองโกเลีย ปากีสถาน ตุรกี ฯลฯ ลักษณะพืชพันธุ์ที่สำคัญคือ^{ลักษณะ} ความเข้มแข็ง ทนทาน สามารถเจริญ^{เจริญดี} บนดินทรายและดินหิน ต้านทานต่อ^{ต้านทาน} โรคเชื้อราและแมลงศรีษะ แต่ต้องระวัง^{ระวัง} ไม่ให้โดนน้ำท่วมนานๆ 以免^{以免} ทำให้เสียหาย สำหรับการเก็บเกี่ยว ควรเลือกเมล็ดที่^{เมล็ด} ขนาดใหญ่และสุก成熟的 เนื่องจากเมล็ดที่เล็กและไม่สุก อาจไม่สามารถเจริญเติบโต^{เติบโต} ได้ ผลลัพธ์ที่ได้มาจะเป็นแอปเปิลที่อร่อยและมีคุณภาพดี สามารถนำมาทำ成^{ทำ成} อาหาร เช่น น้ำเชื่อม น้ำผลไม้ หรือแม้แต่เป็นเครื่องดื่ม เช่น ชา กาแฟ ฯลฯ ที่มีประโยชน์ต่อสุขภาพ ดังนั้น การปลูกแอปเปิลจึงเป็นหนึ่งในอาชีพเกษตรกรรมที่สำคัญและมีมูลค่าทางเศรษฐกิจ ของประเทศไทย

Ground truth text

experiment 1:
benchmarking on
wiki. screenshots

OCR performance on Wikipedia screenshots

	Character Error Rate			Word Error Rate		
	EasyOCR	Tesseract	GOT	EasyOCR	Tesseract	GOT
English	0.21	0.21	0.67	0.27	0.29	0.67
Indonesian	0.28	0.28	0.61	0.36	0.42	0.71
Vietnamese	0.47	0.38	-	0.45	0.39	-
Thai	0.35	0.55	-	1.45	1.45	-

* Lower → Better

OCR performance on Wikipedia screenshots

	Character Error Rate			Word Error Rate		
	EasyOCR	Tesseract	GOT	EasyOCR	Tesseract	GOT
English	0.21	0.21	0.67	0.27	0.29	0.67
Indonesian	0.28	0.28	0.61	0.36	0.42	0.71
Vietnamese	0.47	0.38	-	0.45	0.39	-
Thai	0.35	0.55	-	1.45	1.45	-

* Lower → Better

- Error rate: **Latin** script

OCR performance on Wikipedia screenshots

	Character Error Rate			Word Error Rate		
	EasyOCR	Tesseract	GOT	EasyOCR	Tesseract	GOT
English	0.21	0.21	0.67	0.27	0.29	0.67
Indonesian	0.28	0.28	0.61	0.36	0.42	0.71
Vietnamese	0.47	0.38	-	0.45	0.39	-
Thai	0.35	0.55	-	1.45	1.45	-

* Lower → Better

- Error rate: **Latin** script < Latin script with **diacritics**

OCR performance on Wikipedia screenshots

	Character Error Rate			Word Error Rate		
	EasyOCR	Tesseract	GOT	EasyOCR	Tesseract	GOT
English	0.21	0.21	0.67	0.27	0.29	0.67
Indonesian	0.28	0.28	0.61	0.36	0.42	0.71
Vietnamese	0.47	0.38	-	0.45	0.39	-
Thai	0.35	0.55	-	1.45	1.45	-

* Lower → Better

- Error rate: **Latin** script < Latin script with **diacritics** < **Brahmic** script

OCR performance on Wikipedia screenshots

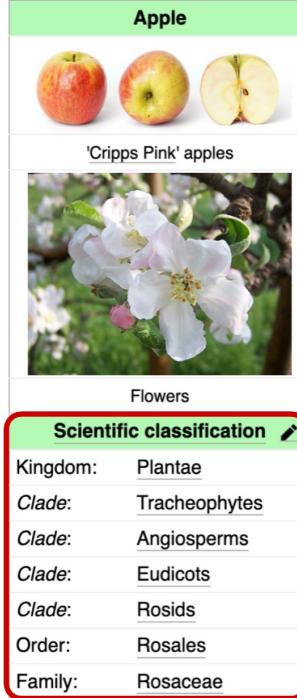
	Character Error Rate			Word Error Rate		
	EasyOCR	Tesseract	GOT	EasyOCR	Tesseract	GOT
English	0.21	0.21	0.67	0.27	0.29	0.67
Indonesian	0.28	0.28	0.61	0.36	0.42	0.71
Vietnamese	0.47	0.38	-	0.45	0.39	-
Thai	0.35	0.55	-	1.45	1.45	-
Average	0.33	0.35	0.64			* Lower → Better

- Error rate: Latin script < Latin script with diacritics < Brahmic script
- **EasyOCR** achieved the best overall performance

Limitation: Complex layouts and multi-modal elements **interfere** with OCR results.

An **apple** is a round, edible fruit produced by an **apple tree** (*Malus spp.*, among them the **domestic or orchard apple**; *Malus domestica*). Apple trees are cultivated worldwide and are the most widely grown species in the genus *Malus*. The tree originated in Central Asia, where its wild ancestor, *Malus sieversii*, is still found. Apples have been grown for thousands of years in Eurasia and were introduced to North America by European colonists. Apples have religious and mythological significance in many cultures, including Norse, Greek, and European Christian tradition.

Apples grown from seed tend to be very different from those of their parents, and the resultant fruit frequently lacks desired characteristics. For commercial purposes, including botanical evaluation, apple cultivars are propagated by clonal grafting onto rootstocks. Apple trees grown without rootstocks tend to be larger and much slower to fruit after planting. Rootstocks are used to control the speed of growth and the size of the resulting tree, allowing for



Apples grown from seed tend to be very different from those of their parents, and the resultant fruit frequently lacks desired characteristics. For commercial purposes, including botanical evaluation, apple cultivars are propagated by clonal grafting onto rootstocks. Apple trees grown without rootstocks tend to be larger and much slower to fruit after planting. Rootstocks are used to control the speed of growth and the size of the resulting tree, allowing for



Inflated CERs and WERs for Wikipedia screenshots

Need for **layout analysis** to process digital docs.

01

Motivation

02

Related Work

03

Experiment 1: Benchmarking on
Real-world Wikipedia Screenshots

04

Experiment 2: Benchmarking on
Synthetic Data

05

Experiment 3: Fine-tuning

Synthetic data allows **controlled distortions** to dataset
and **minimizes** annotation errors



- Goals:
1. Analyze OCR robustness against **noise** on SEA languages (RQ1)
 2. **Isolate** script-related errors (RQ2)

1. Analyze OCR robustness against **noise**

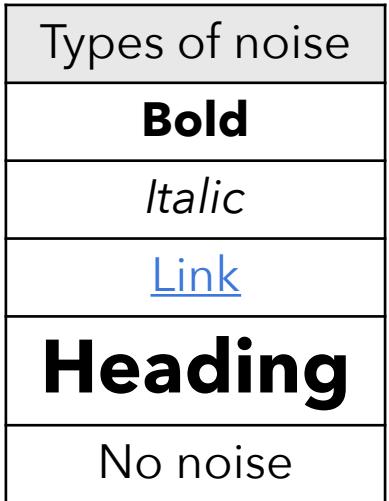
Types of noise
Bold
<i>Italic</i>
<u>Link</u>
Heading
No noise

Why?

1. **Commonly** found in digital text and web-based content
2. Each noise type presents a **unique challenge** for OCR systems

experiment 2: benchmarking on synthetic data

1. Analyze OCR robustness against **noise**



Generate document images with noise using WeasyPrint



Run and evaluate OCR tools on each dataset

experiment 2:
benchmarking on
synthetic data

OCR performance on noise-free synthetic data

	Character Error Rate		
	EasyOCR	Tesseract	GOT
English	0.03	0.02	0.01
Indonesian	0.02	0.02	0.01
Vietnamese	0.10	0.03	-
Thai	0.07	0.09	-

* Lower → Better

Similar trend
as benchmark
on Wiki.
screenshots

OCR performance on noise-free synthetic data

	Character Error Rate		
	EasyOCR	Tesseract	GOT
English	0.03	0.02	0.01
Indonesian	0.02	0.02	0.01
Vietnamese	0.10	0.03	-
Thai	0.07	0.09	-

Average

0.02

0.02

0.06

0.08

* Lower → Better

- Error rate: **Latin** script < Latin script with **diacritics** < **Brahmic** script

OCR performance on noise-free synthetic data

	Character Error Rate		
	EasyOCR	Tesseract	GOT
English	0.03 ↓	0.02 ↓	0.01 ↓
Indonesian	0.02 ↓	0.02 ↓	0.01 ↓
Vietnamese	0.10 ↓	0.03 ↓	-
Thai	0.07 ↓	0.09 ↓	-

* Lower → Better

Avg. decrease
of **0.34**

- Error rate: Latin script < Latin script with diacritics < Brahmic script
- OCR tools performed **better** on synthetic data than on Wiki. screenshots

OCR performance on noise-free synthetic data

	Character Error Rate		
	EasyOCR	Tesseract	GOT
English	0.03	0.02	0.01
Indonesian	0.02	0.02	0.01
Vietnamese	0.10	0.03	-
Thai	0.07	0.09	-

Average for
synthetic data 0.05 > 0.04 > 0.01

Average for Wiki.
screenshots 0.33 < 0.35 < 0.64

Tesseract and
GOT potentially
**less robust to
noise**

- Error rate: Latin script < Latin script with diacritics < Brahmic script
- OCR tools performed better on synthetic data than on Wiki. screenshots
- Tesseract and GOT **performed well** on synthetic data, but poorly on Wikipedia screenshots

Impact of noise on OCR accuracy

EasyOCR

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.03	-7.1%	0.0%	3.6%	-7.1%
Indonesian	0.02	-5.6%	0.0%	22.2%	-5.6%
Vietnamese	0.10	-1.0%	-1.9%	-6.7%	-11.5%
Thai	0.07	0.0%	0.0%	7.2%	0.0%

Tesseract

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.02	-4.2%	0.0%	0.0%	8.3%
Indonesian	0.02	0.0%	0.0%	0.0%	21.7%
Vietnamese	0.03	0.0%	3.8%	3.8%	26.9%
Thai	0.09	3.2%	0.0%	5.4%	4.3%

GOT

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.01	27.3%	36.4%	-18.2%	218.2%
Indonesian	0.01	9.1%	18.2%	36.4%	63.6%

Impact of noise on OCR accuracy

EasyOCR

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.03	-7.1%	0.0%	3.6%	-7.1%
Indonesian	0.02	-5.6%	0.0%	22.2%	-5.6%
Vietnamese	0.10	-1.0%	-1.9%	-6.7%	-11.5%
Thai	0.07	0.0%	0.0%	7.2%	0.0%

Tesseract

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.02	-4.2%	0.0%	0.0%	8.3%
Indonesian	0.02	0.0%	0.0%	0.0%	21.7%
Vietnamese	0.03	0.0%	3.8%	3.8%	26.9%
Thai	0.09	3.2%	0.0%	5.4%	4.3%

GOT

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.01	27.3%	36.4%	-18.2%	218.2%
Indonesian	0.01	9.1%	18.2%	36.4%	63.6%

- Tesseract and GOT
sensitive to noise

Impact of noise on OCR accuracy

EasyOCR

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.03	-7.1%	0.0%	3.6%	-7.1%
Indonesian	0.02	-5.6%	0.0%	22.2%	-5.6%
Vietnamese	0.10	-1.0%	-1.9%	-6.7%	-11.5%
Thai	0.07	0.0%	0.0%	7.2%	0.0%

Tesseract

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.02	-4.2%	0.0%	0.0%	8.3%
Indonesian	0.02	0.0%	0.0%	0.0%	21.7%
Vietnamese	0.03	0.0%	3.8%	3.8%	26.9%
Thai	0.09	3.2%	0.0%	5.4%	4.3%

GOT

	No Noise CER	Bold % Change	Italic % Change	Link % Change	Heading % Change
English	0.01	27.3%	36.4%	-18.2%	218.2%
Indonesian	0.01	9.1%	18.2%	36.4%	63.6%

Average

2.2%

5.6%

5.4%

31.9%

- Tesseract and GOT sensitive to noise
- Heading noise presents greatest challenge for OCR tools

2. **Classify** errors by character type

	Included Characters
Arabic digit	0-9
Latin letter	a-z, A-Z
Latin letter with diacritic	à-ÿ, À-Ŷ
Vietnamese special letter	đ, Đ
Thai letter	ກ-ຂ
Special symbol	.,!?:();"-'--\$%/&+-=[]{}
Whitespace	—
Other	

OCR results from synthetic dataset **without noise**



Identify and **classify** misrecognized characters using RegEx

Error classification in English and Indonesian articles

English

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	38,324	0.7%	1.9%	0.3%
Latin letter	1,546,964	1.3%	1.8%	0.4%
Latin letter with diacritic	452	100.0%	56.0%	15.9%
Special symbol	55,110	29.5%	3.7%	4.3%
Whitespace	317,587	4.9%	4.3%	3.6%
Other	1,562	103.6%	92.6%	127.9%

Articles with Latin scripts:

Indonesian

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	24,947	0.4%	1.8%	0.2%
Latin letter	1,208,707	0.5%	1.8%	0.4%
Latin letter with diacritic	276	7.2%	100.0%	16.7%
Special symbol	38,980	22.9%	4.9%	1.4%
Whitespace	207,556	4.8%	5.1%	4.1%
Other	1,262	94.4%	99.6%	64.8%

Error classification in English and Indonesian articles

English

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	38,324	0.7%	1.9%	0.3%
Latin letter	1,546,964	1.3%	1.8%	0.4%
Latin letter with diacritic	452	100.0%	56.0%	15.9%
Special symbol	55,110	29.5%	3.7%	4.3%
Whitespace	317,587	4.9%	4.3%	3.6%
Other	1,562	103.6%	92.6%	127.9%

Indonesian

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	24,947	0.4%	1.8%	0.2%
Latin letter	1,208,707	0.5%	1.8%	0.4%
Latin letter with diacritic	276	7.2%	100.0%	16.7%
Special symbol	38,980	22.9%	4.9%	1.4%
Whitespace	207,556	4.8%	5.1%	4.1%
Other	1,262	94.4%	99.6%	64.8%

Articles with Latin scripts:

- High accuracies on **Latin letters** and **Arabic digits**

Error classification in English and Indonesian articles

English

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	38,324	0.7%	1.9%	0.3%
Latin letter	1,546,964	1.3%	1.8%	0.4%
Latin letter with diacritic	452	100.0%	56.0%	15.9%
Special symbol	55,110	29.5%	3.7%	4.3%
Whitespace	317,587	4.9%	4.3%	3.6%
Other	1,562	103.6%	92.6%	127.9%

Indonesian

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	24,947	0.4%	1.8%	0.2%
Latin letter	1,208,707	0.5%	1.8%	0.4%
Latin letter with diacritic	276	7.2%	100.0%	16.7%
Special symbol	38,980	22.9%	4.9%	1.4%
Whitespace	207,556	4.8%	5.1%	4.1%
Other	1,262	94.4%	99.6%	64.8%

Articles with Latin scripts:

- High accuracies on **Latin letters** and **Arabic digits**
- Lower accuracies for **Latin letters with diacritics**

Error classification in English and Indonesian articles

English

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	38,324	0.7%	1.9%	0.3%
Latin letter	1,546,964	1.3%	1.8%	0.4%
Latin letter with diacritic	452	100.0%	56.0%	15.9%
Special symbol	55,110	29.5%	3.7%	4.3%
Whitespace	317,587	4.9%	4.3%	3.6%
Other	1,562	103.6%	92.6%	127.9%

Indonesian

	Count	EasyOCR % Missed	Tesseract % Missed	GOT % Missed
Arabic digit	24,947	0.4%	1.8%	0.2%
Latin letter	1,208,707	0.5%	1.8%	0.4%
Latin letter with diacritic	276	7.2%	100.0%	16.7%
Special symbol	38,980	22.9%	4.9%	1.4%
Whitespace	207,556	4.8%	5.1%	4.1%
Other	1,262	94.4%	99.6%	64.8%

Articles with Latin scripts:

- High accuracies on **Latin letters** and **Arabic digits**
- Lower accuracies for Latin letters with **diacritics**
- Most **uncategorized characters** misclassified
 - E.g., Greek, Cyrillic, Arabic, Chinese characters

Error classification in Vietnamese and Thai articles

Vietnamese

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	31,473	1.1%	2.2%
Latin letter	916,667	8.5%	1.8%
Latin letter with diacritic	294,406	14.9%	1.9%
Vietnamese special letter	30,903	4.2%	1.8%
Special symbol	41,655	25.4%	3.1%
Whitespace	367,936	10.9%	5.3%
Other	1,909	105.4%	88.6%

Thai

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	22,580	0.9%	6.7%
Latin letter	36,174	100.0%	100.0%
Latin letter with diacritic	101	100.0%	100.0%
Thai letter	1,014,395	1.2%	4.2%
Special symbol	14,246	8.0%	10.6%
Whitespace	58,164	37.5%	37.2%
Other	890	105.2%	110.4%

Articles with more complex scripts:

Error classification in Vietnamese and Thai articles

Vietnamese

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	31,473	1.1%	2.2%
Latin letter	916,667	8.5%	1.8%
Latin letter with diacritic	294,406	14.9%	1.9%
Vietnamese special letter	30,903	4.2%	1.8%
Special symbol	41,655	25.4%	3.1%
Whitespace	367,936	10.9%	5.3%
Other	1,909	105.4%	88.6%

Thai

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	22,580	0.9%	6.7%
Latin letter	36,174	100.0%	100.0%
Latin letter with diacritic	101	100.0%	100.0%
Thai letter	1,014,395	1.2%	4.2%
Special symbol	14,246	8.0%	10.6%
Whitespace	58,164	37.5%	37.2%
Other	890	105.2%	110.4%

Articles with more complex scripts:

- Tesseract achieved high accuracies on **Vietnamese letters**

Error classification in Vietnamese and Thai articles

Vietnamese

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	31,473	1.1%	2.2%
Latin letter	916,667	8.5%	1.8%
Latin letter with diacritic	294,406	14.9%	1.9%
Vietnamese special letter	30,903	4.2%	1.8%
Special symbol	41,655	25.4%	3.1%
Whitespace	367,936	10.9%	5.3%
Other	1,909	105.4%	88.6%

Thai

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	22,580	0.9%	6.7%
Latin letter	36,174	100.0%	100.0%
Latin letter with diacritic	101	100.0%	100.0%
Thai letter	1,014,395	1.2%	4.2%
Special symbol	14,246	8.0%	10.6%
Whitespace	58,164	37.5%	37.2%
Other	890	105.2%	110.4%

Articles with more complex scripts:

- Tesseract achieved high accuracies on **Vietnamese letters**
- High accuracies on **Thai letters**

Error classification in Vietnamese and Thai articles

Vietnamese

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	31,473	1.1%	2.2%
Latin letter	916,667	8.5%	1.8%
Latin letter with diacritic	294,406	14.9%	1.9%
Vietnamese special letter	30,903	4.2%	1.8%
Special symbol	41,655	25.4%	3.1%
Whitespace	367,936	10.9%	5.3%
Other	1,909	105.4%	88.6%

Thai

	Count	EasyOCR % Missed	Tesseract % Missed
Arabic digit	22,580	0.9%	6.7%
Latin letter	36,174	100.0%	100.0%
Latin letter with diacritic	101	100.0%	100.0%
Thai letter	1,014,395	1.2%	4.2%
Special symbol	14,246	8.0%	10.6%
Whitespace	58,164	37.5%	37.2%
Other	890	105.2%	110.4%

Articles with more complex scripts:

- Tesseract achieved high accuracies on **Vietnamese letters**
- High accuracies on **Thai letters**
- Thai OCR models did not recognize **Latin letters with or without diacritics**

Implication 1: OCR tools recognize **Viet.** and **Thai scripts** well but struggle with **non-native characters**.

- Related work:
 - OCR tools perform **poorly** on complex scripts¹
 - Latin scripts perform **better** than SEA scripts²
- Our results: **Latin** script < Latin script with **diacritics** < **Brahmic** script
 - High accuracies on **Vietnamese letters**
 - High accuracies on **Thai letters**
 - Thai OCR models did not recognize **Latin letters with or without diacritics**

¹ A concise survey of OCR for low-resource languages (Agarwal & Anastasopoulos, 2024)

² OCR improves machine translation for low-resource languages (Ignat et al., 2022)

Implication 1: OCR tools recognize **Viet.** and **Thai scripts** well but struggle with **non-native characters**.

- Related work:
 - OCR tools perform **poorly** on complex scripts¹
 - Latin scripts perform **better** than SEA scripts²
 - Our results: **Latin** script < Latin script with **diacritics** < **Brahmic** script
 - High accuracies on **Vietnamese letters**
 - High accuracies on **Thai letters**
 - Thai OCR models did not recognize **Latin letters with or without diacritics**
- 
- Challenges** common belief

¹ A concise survey of OCR for low-resource languages (Agarwal & Anastasopoulos, 2024)

² OCR improves machine translation for low-resource languages (Ignat et al., 2022)

Implication 1: OCR tools recognize **Viet.** and **Thai scripts** well but struggle with **non-native characters**.

Misclassification of characters
outside of target language
inflates overall CER

- High accuracies on **Vietnamese letters**
- High accuracies on **Thai letters**
- Thai OCR models did not recognize
Latin letters with or without diacritics

OCR tools struggle with
multilingual content

Challenges common belief

Need for OCR systems to better handle **mixed-script content**

Goal: Explore how fine-tuning can enhance OCR accuracy on SEA languages (RQ3)



Fine-tune GOT on Vietnamese and Thai

Why?

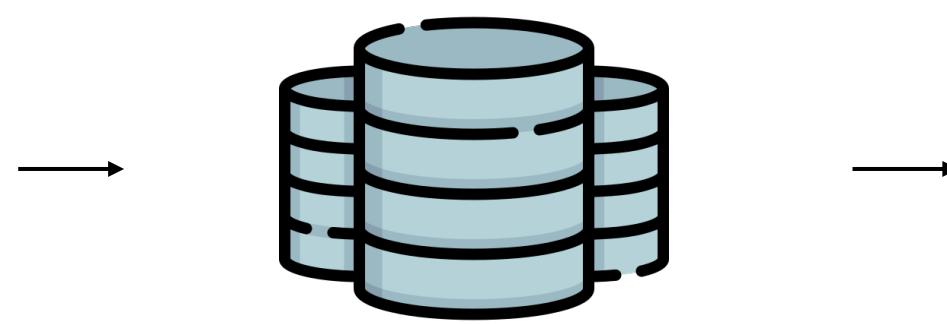
1. GOT does **not** officially support Vietnamese and Thai
2. Related work: Fine-tuning pre-trained **transformer-based** models is efficient in adapting new languages
3. GOT already achieved state-of-the-art results on **Indonesian**

Fine-tune GOT on Vietnamese and Thai



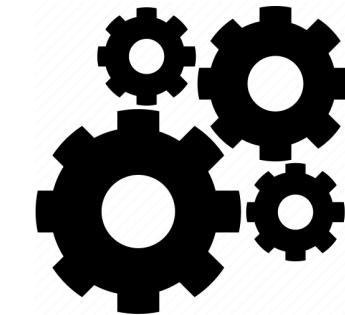
Generate more noise-free synthetic data

960 training images
+ **50** test images



Split dataset in sizes
of increasing scale
to approximate diminishing returns

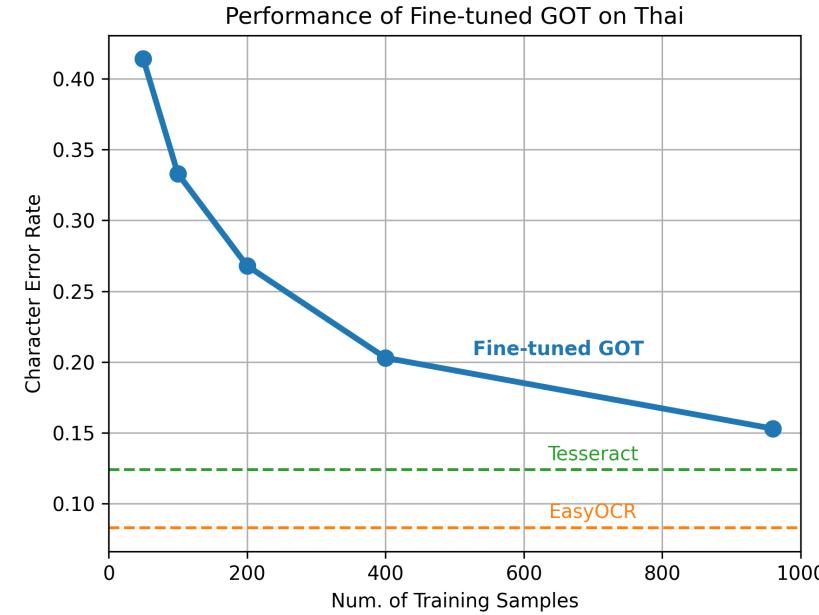
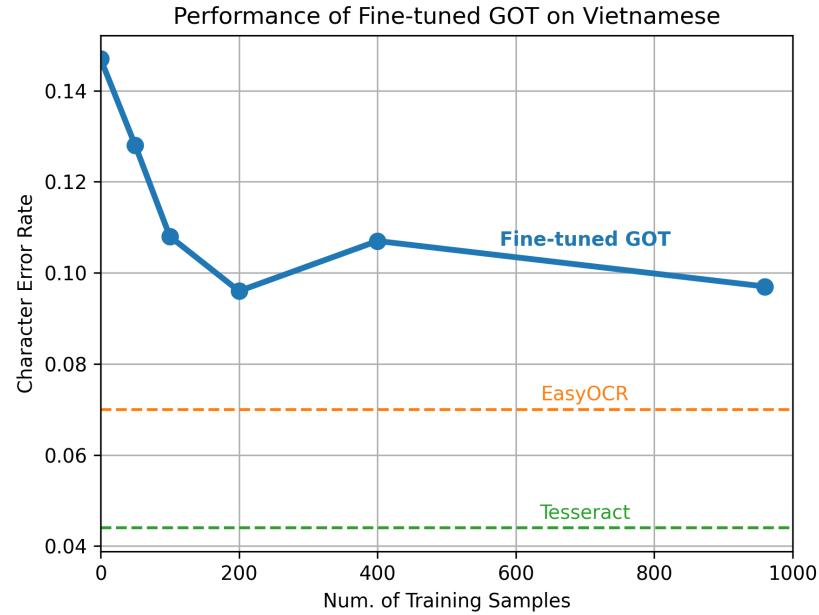
Datasets with 50, 100, 200, 400, 960 images



Fine-tune GOT using each dataset

Fine-tuned GOTs

experiment 3: fine-tuning



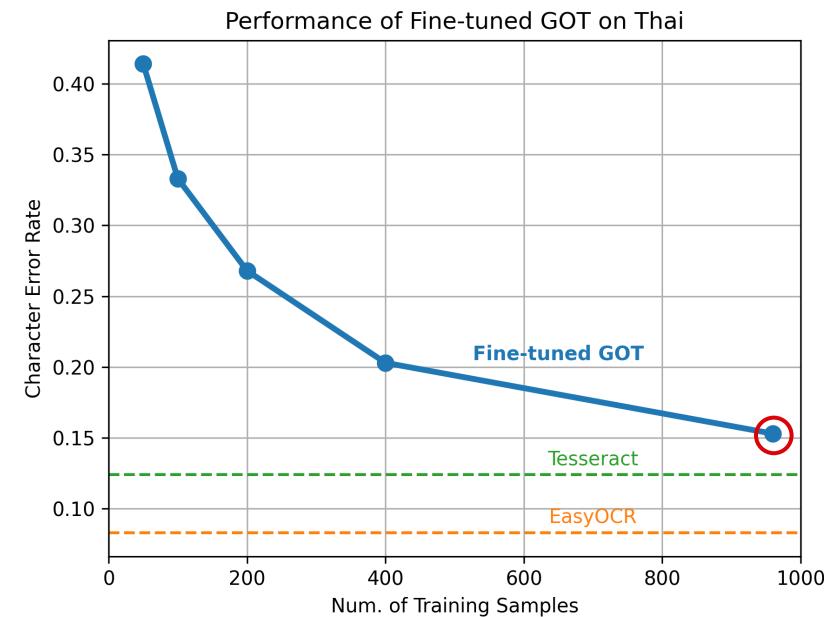
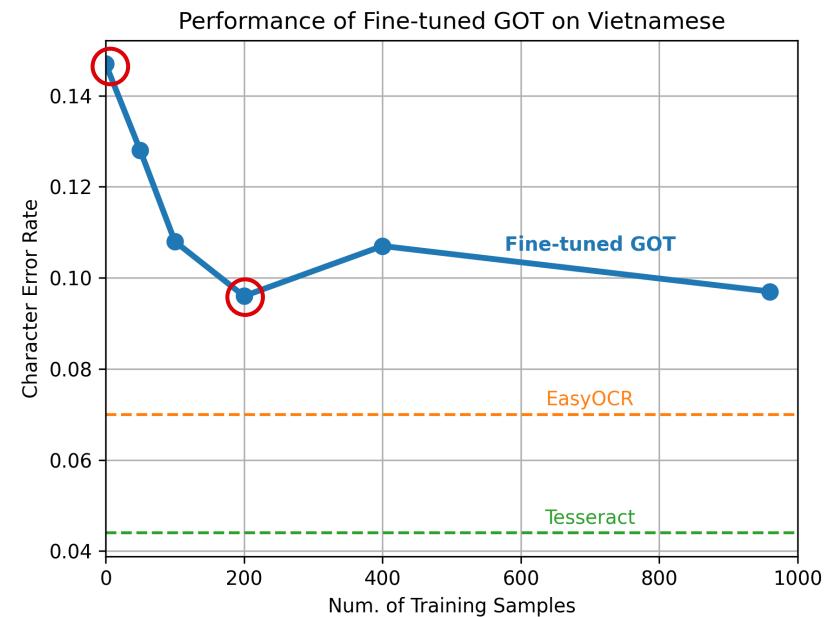
experiment 3: fine-tuning

No fine-tuning:

-33%

200 samples:

0.10



No fine-tuning:

Cannot recognize

960 samples:

0.15

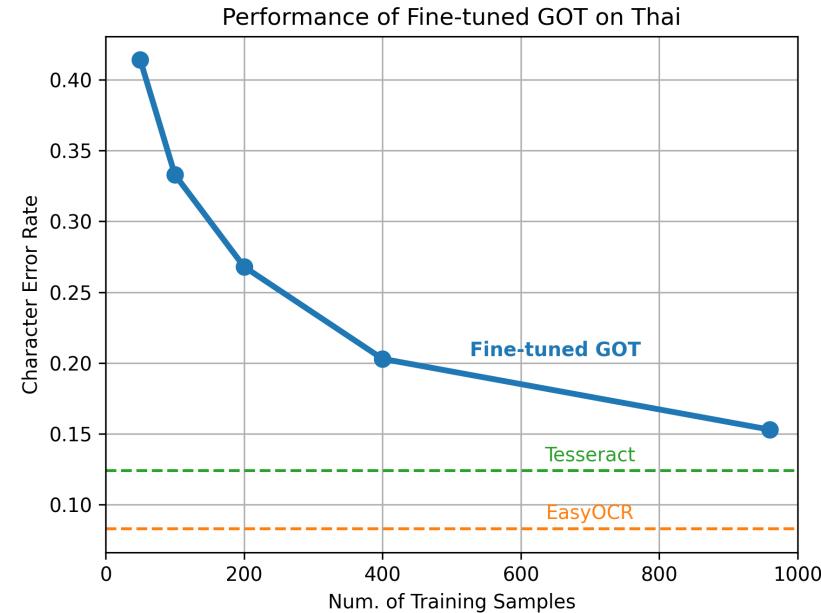
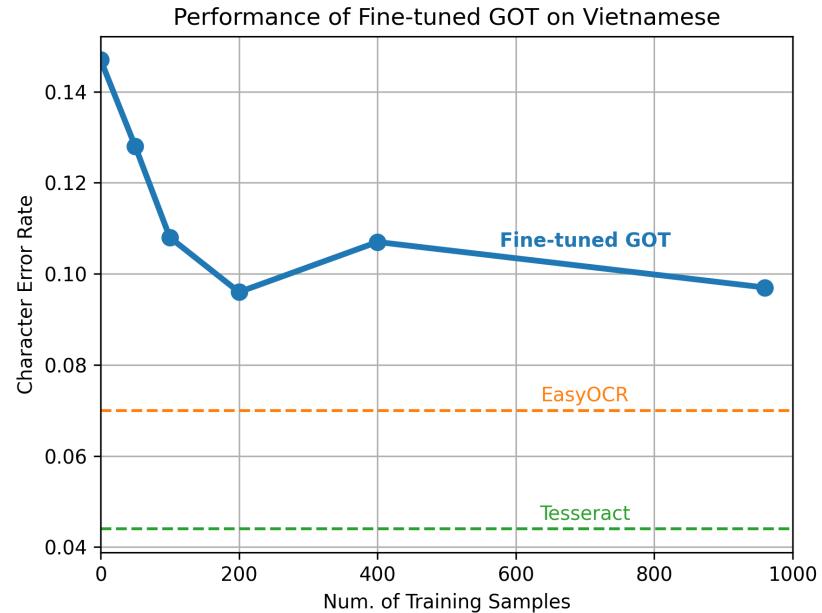
- Fine-tuning **improved** performance with relatively little data

experiment 3: fine-tuning

Ours: **0.10**

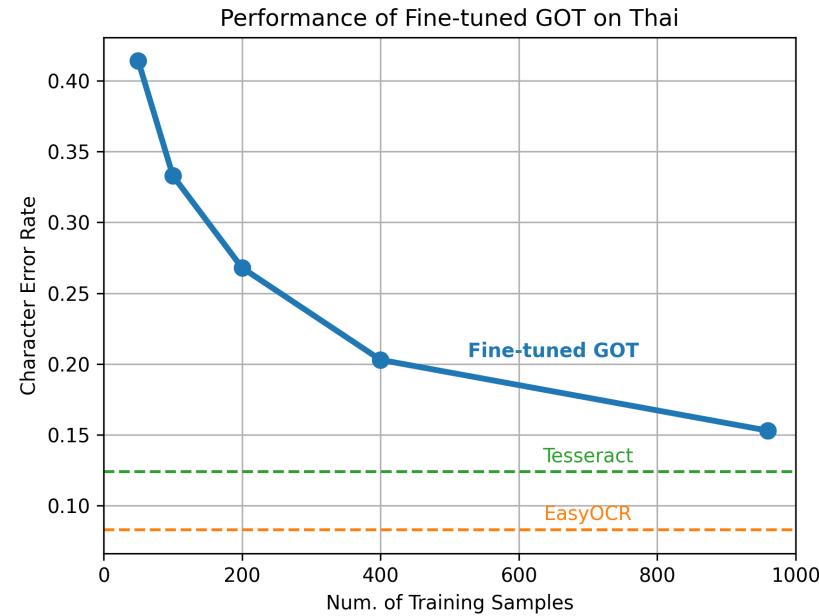
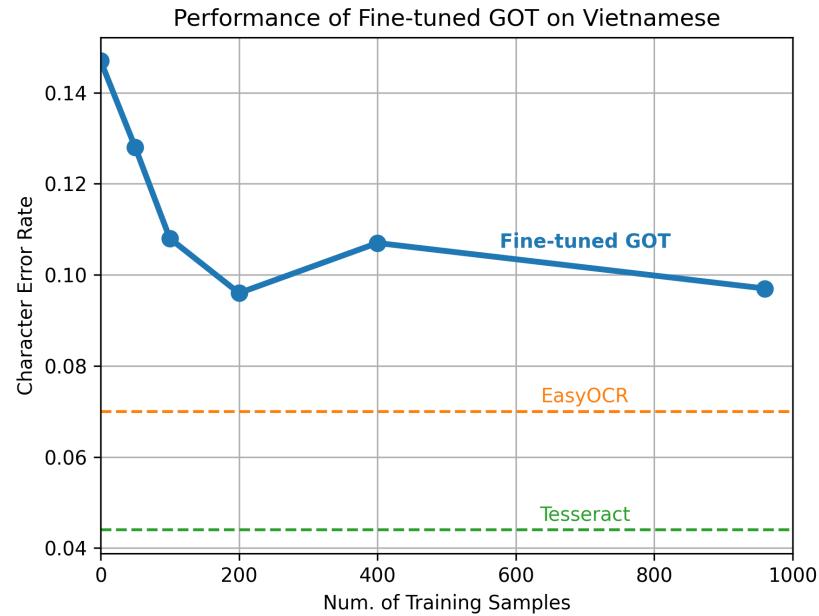
EasyOCR: 0.07

Tesseract: 0.04



Ours: **0.15**
Tesseract: 0.12
EasyOCR: 0.08

- Fine-tuning **improved** performance with relatively little data
- Fine-tuned GOT **underperformed** relative to EasyOCR and Tesseract



- Fine-tuning **improved** performance with relatively little data
- Fine-tuned GOT **underperformed** relative to EasyOCR and Tesseract
- Model performance did **not** consistently improve with more data

Implication 2: Fine-tuning is **effective** in adapting OCR systems to new SEA languages.

Fine-tuning **improved** performance with relatively little data



Practical approach for adapting to new **low-resource** SEA languages

Surpassing state-of-the-art OCR tools may require exploring **alternative models** and optimizing **fine-tuning strategies**

 Goal: **Benchmark** and **improve** OCR for SEA languages

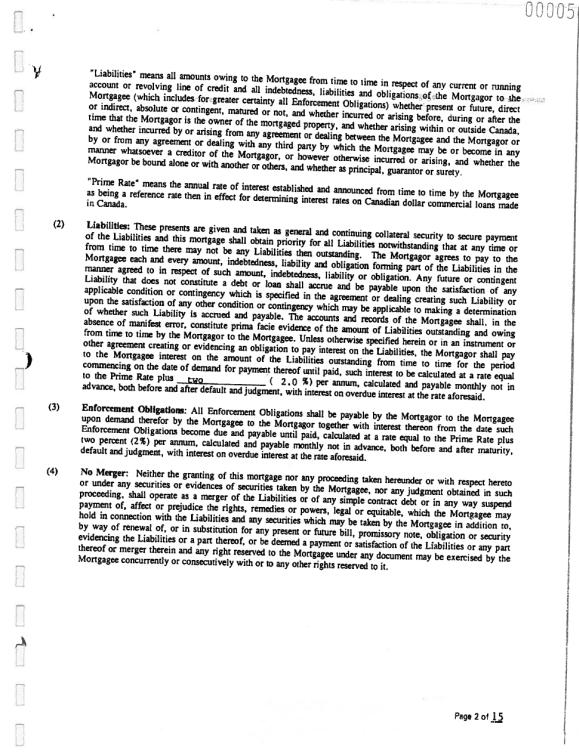
What I did:

- 1 **Benchmarked** EasyOCR, Tesseract, and GOT on English, Indo., Viet., and Thai
- 2 Created **reusable pipeline** for collecting Wiki. data and benchmarking OCR tools
- 3 **Fine-tuned** GOT on Vietnamese and Thai

Key findings & implications:

- ✨ OCR tools recognize **Vietnamese** and **Thai scripts** well, but struggle with non-native characters → Need for **mixed-script OCR**
- ✨ Fine-tuning is **effective** in adapting OCR systems to new SEA languages

Limitation: Wikipedia articles **cannot fully replace** scanned physical documents.



WIKIPEDIA
The Free Encyclopedia

Apple

An apple is a round, edible fruit produced by an **apple tree** (*Malus spp.*, among them the **domestic** or **orchard apple**; *Malus domestica*). Apple trees are cultivated worldwide and are the most widely grown species in the genus *Malus*. The tree originated in Central Asia, where its wild ancestor, *Malus sieversii*, is still found. Apples have been grown for thousands of years in Eurasia and were introduced to North America by European colonists. Apples have religious and mythological significance in many cultures, including Norse, Greek, and European Christian tradition.

Apples grown from seed tend to be very different from those of their parents, and the resultant fruit frequently lacks desired characteristics. For commercial purposes, including botanical evaluation, apple cultivars are propagated by clonal grafting onto rootstocks. Apple trees grown without rootstocks tend to be larger and much slower to fruit after planting. Rootstocks are used to control the speed of growth and the size of the resulting tree, allowing for easier harvesting.

There are more than 7,500 cultivars of apples. Different cultivars are bred for various tastes and uses, including cooking, eating raw, and cider or apple juice production. Trees and fruit are prone to fungal, bacterial, and pest problems, which can be controlled by a number of organic and non-organic means. In 2010, the fruit's genome was sequenced as part of research on disease control and selective breeding in apple production.

Scientific classification

Kingdom:	Plantae
Clade:	Tracheophytes
Clade:	Angiosperms
Clade:	Eudicots
Clade:	Rosids
Order:	Rosales
Family:	Rosaceae
Genus:	<i>Malus</i>
Species:	<i>M. domestica</i>
Binomial name	
<i>Malus domestica</i> (Suckow) Borkh., 1803	
Synonyms ^{[1][2]}	
• <i>M. communis</i> Desf., 1768	
• <i>M. pumila</i> Mill.	
• <i>M. frutescens</i> Medik.	

Physical documents have **more diverse noise** types and **practical use cases**

Why still use Wikipedia?

1. Similar **linguistic** features
2. **More accessible** to obtain lots of data
3. Growing use cases of OCR on **digital screenshots**