# សាកលវិទ្យាល័យភូមិន្ទភ្នំពេញ
## ROYAL UNIVERSITY OF PHNOM PENH

ការស្រាវជ្រាវវិធីសាស្ត្រថ្មី សម្រាប់កំណត់សម្គាល់អត្ថបទអក្សរខ្មែរទូទៅ និង
បានប្រើប្រាស់ស្ថាបត្យកម្ម Craft ជាមួយនឹង TrOCR

# A novel End-to-End approach for General Khmer Text Recognition using Craft with TrOCR Architecture

Mr. Vitou Soy

## A Thesis

In Partial Fulfilment of the Requirement for the Degree of
Bachelor of Engineering in Information-Technology-Engineering

Examination committee:   Mr. Sokchea Kor (Advisor)
Mr. Chanpiseth Chap (committee)
Mrs. Daly Chea (committee)
Dr. ........................

June 2025

# មូលន័យសង្ខេប

ក្នុងសហគមន៍បច្ចេកវិទ្យាព័ត៌មានសម័យថ្មី ការចាប់យកអក្ខរបទចេញពីរូបភាព – [ OCR ] ( Optical Character Recognition ) ក្លាយជាបច្ចេកវិទ្យាសំខាន់មួយដែលត្រូវបានប្រើប្រាស់ យ៉ាងទូលំទូលាយ សម្រាប់បំលែងឯកសារ សរសេរ ឬរូបភាពអក្សរឲ្យទៅជាអក្ខរទ អេឡិចត្រូនិច ( digital text )។ ការអភិវឌ្ឍ OCR សម្រាប់ភាសាខ្មែរ តែងតែ ប្រឈមនឹងបញ្ហាជាច្រើន ដោយសារកង្វះនៃប្រភពទិន្នន័យ និងឯកសារសម្រាប់ train AI model។ ដើម្បីដោះស្រាយ បញ្ហានេះ យើងបានបង្កើតទិន្នន័យសិប្បនិម្មិត ( Synthetic Dataset ) ដោយប្រើវិធីសាស្ត្របច្ចេកទេសកម្រិតខ្ពស់។

ក្នុងដំណើរការបង្កើតទិន្នន័យសិប្បនិម្មិត ( Synthetic Dataset ) រួមមាន៖

- វិធីសាស្ត្រក្នុងការប្រមូលអក្ខរបទចេញពីអ៊ីនធឺណិត មានដូចខាងក្រោម ( Scrape data ) ៖
  - ដំណាក់កាលទីមួយ៖ យើងបានប្រមូលអក្ខរបទចេញពី khsearch.com, Chuon-Nath-Dictionary, Alpha-Word, Google-Word, និងចុងក្រោយគឺ Huggingface.com ។
  - ដំណាក់កាលទីពីរ៖ យើងបានសម្អាត ទិន្នន័យទាំងអស់នោះ ឆ្លងកាត់ដំណើរការ ដូចជា លុបចោលគ្មាន អក្សរណាដែលមិនស្រួមាន វត្តមាននៅលើ រូបភាព ញឹកញាប់ និងបានលុបចោល គ្មានអក្សរណាដែល Fonts renders អត់ចេញ។
  - ដំណាក់កាលទីបី៖ ដំណាក់កាលមួយនេះ យើងបានធ្វើការ កាត់ប្រយោគទាំងអស់នោះ ជាពាក្យៗ ដោយ ប្រើប្រាស់ library ឈ្មោះ khmer-nltk
  - ដំណាក់កាលទីបួន៖ ចុងក្រោយ ក៏បានរៀបចំជា ប្រយោគដែល មានប្រវែង Random ពី ១ អក្សរ រហូតដល់ ១១០ អក្សរ ។

- បង្កើតរូបភាពដោយអនុវត្តតាមលក្ខខណ្ឌខាងក្រោម ៖
  - ផ្ទៃខាងក្រោយចែជនុ ( Apply Different backgrounds )
  - បំពាក់ពុម្ពអក្សរផ្សេងៗគ្នា ( Apply Different fonts )
  - Noise: `gaussian_noise`, `salt_pepper_noise`, `speckle_noise`, `blur`
  - បង្វិលអក្សរបន្តិច ( random rotation text )
  - បញ្ចូល Margin Randomly ( 1, 5 ) pixels

- សរុបមកយើងបានបង្កើត Data ជាង ៤ លាន records សម្រាប់ train OCR model

Architecture OCR ត្រូវបានបែងចែកជា ២ ផ្នែក៖ Text Detection និង Text Recognition:

Text Detection: យើងប្រើម៉ូដែល CRAFT ដោយបានធ្វើការ Train ឡើងវិញដោយ បាន annotation ទៅលើ លើរូបភាពប្រហែល ៥០០ images និងសរុបចំនួន bounding box ជាង ១០,០០០ boxes។

Text Recognition: យើងប្រើ TrOCR base model ចេញពី Microsoft ( មាននៅក្នុង Hugging Face ) ហើយ បាន fine-tune ទៅលើ dataset ខ្មែរសិប្បនិម្មិត ( Synthetic Dataset ) ដើម្បីបង្កើន សមត្ថភាពក្នុង ការសម្គាល់ អក្សរខ្មែរ។

លទ្ធផលសិក្សាបានបង្ហាញថា OCR របស់ពួកយើងអាចសម្គាល់អក្ខរបទចេញពីរូបភាព បានដោយភាពត្រឹមត្រូវ លើសពី ៩០%។ ដូច្នេះ ការសិក្សានេះបង្ហាញអំពី សក្តានុពលនៃការបង្កើត dataset និងការប្រើម៉ូដែលដ៏នាន់ថ្មី ដើម្បីអភិវឌ្ឍន៍ OCR ភាសាខ្មែរឲ្យមានប្រសិទ្ធភាពកាន់តែខ្ពស់។ វ៉ាមានសមត្ថភាព អាចចាប់យកអក្ខរបទមិនត្រឹមតែ ពាក្យខ្លីៗ ប៉ុណ្ណោះទេ តែវ៉ាក៏អាចធ្វើការចាប់យក ដូចជា មួយតួអក្សរដោយមួយតួអក្សរ, ពាក្យដោយពាក្យ, ប្រយោគ ដោយប្រយោគ រហូតដល់ មួយប្រយោគវែង ១១០ តួអក្សរថែមទៀតផងផង ។ ហើយលើសពីនោះទៀត វ៉ាក៏អាចធ្វើការ កំណត់សម្គាល់ទៅលើ ពីរ ភាសាចម្បង ទាំងភាសាខ្មែរ និងភាសាអង់គ្លេស ។

# Abstract

In the modern era of information technology, Optical Character Recognition (OCR) has emerged as a crucial technology for converting printed or handwritten text from images into digital form. However, the development of OCR systems for the Khmer language presents significant challenges, primarily due to the lack of large-scale annotated datasets. To address this limitation, we constructed a high-quality synthetic dataset using an advanced data generation pipeline. Our Khmer OCR system consists of two core components:

- Text Collection: We gathered Khmer text data from various online sources, including khsearch.com, Chuon-Nath Dictionary, Alpha-Word, Google-Word, and Hugging Face.

- Data Cleaning: We processed and cleaned the collected text by removing uncommon characters, symbols that are rarely rendered correctly by fonts, and excessive whitespace.

- Text Segmentation: Sentences were tokenized into words using the khmer-nltk library, and then reconstructed into randomized sentence lengths ranging from 1 to 110 characters.

- Image Generation: We rendered text into synthetic images by:
  - Applying random backgrounds and a variety of Khmer fonts
  - Adding diverse noise types such as Gaussian noise, salt-and-pepper noise, speckle noise, and blur
  - Introducing slight random rotations and random margins (1–5 pixels)

- As a result, we generated over 4 million high-quality synthetic image-text pairs to train the OCR model.

Our Khmer OCR system consists of two core components:

- **Text Detection:** We fine-tuned the CRAFT (Character Region Awareness for Text Detection) model using 500 manually annotated images, totaling over 10,000 bounding boxes.

- **Text Recognition:** We fine-tuned Microsoft's TrOCR base model (available on Hugging Face) on our synthetic Khmer dataset to improve its ability to recognize Khmer text.

The evaluation results demonstrate that our system achieves a recognition accuracy exceeding 90%. These findings highlight the effectiveness of combining synthetic data generation with modern transformer-based architectures to significantly advance Khmer OCR capabilities. Notably, the system can accurately recognize a wide range of text—from single characters and individual words to full sentences of up to 110 characters—and supports both Khmer and English languages.

# SUPERVISOR's RESEARCH SUPERVISION STATEMENT

Name of program: Khmer Studies
Name of candidate: Vitou Soy

Title of research report: A novel End-to-End approach for General Khmer Text Recognition using Craft with TrOCR Architecture

This is to certify that the research carried out for the above titled master's research report was completed by the above named candidate under my direct supervision. This thesis material has not been used for any other degree. The candidate has demonstrated strong research capabilities and independence in developing novel approaches for Khmer text recognition. The research methodology, implementation, and results are original contributions to the field of Khmer OCR technology. I have provided guidance and oversight throughout the research process while allowing the candidate to explore innovative solutions.

Supervisor's name: Sokchea Kor
Supervisor's signature:…………………
Date…………………………………………

# CANDIDATE'S STATEMENT

TO WHOM IT MAY CONCERN

This is to certify that the dissertation that I, Vitou Soy, hereby present, entitled "Advancing Khmer Optical Character Recognition: A Synthetic Data-Driven Approach," for the degree of Bachelor of Engineering in Information Technology at the Royal University of Phnom Penh, is entirely my own work. Furthermore, it has not been used to fulfill the requirements of any other qualification, in the whole or in part, at this or any other University or equivalent institution. The research methodology, implementation, and findings represent original con- tributions to the field of Khmer OCR technology, particularly in developing novel approaches for synthetic data generation and transformer-based text recognition. Through this work, I have demonstrated strong research capabilities and independence in addressing the critical challenges of Khmer text digitization and recognition.

No reference to, or quotation from, this document may be made without the written approval of the author.

Name of Candidate: Vitou Soy

Signed by the candidate:

Date: ...............................................................................

Name of Supervisor: Mr. Sokchea Kor

Countersigned by the Supervisor: ..............................................

Date: ...............................................................................

# ACKNOWLEDGMENTS

# TABLE OF CONTENTS

**Preliminary Pages**

# List of Tables

# List of Figures

# LIST OF ABBREVIATIONS

OCR: Optical Character Recognition . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

CNN: Convolutional Neural Network . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

RNN: Recurrent Neural Network . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

LSTM: Long Short-Term Memory . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

GRU: Gated Recurrent Unit . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

Transformer: Transformer Model . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

BERT: Bidirectional Encoder Representations from Transformers . . . . . . . . . .

TrOCR: Transformer OCR . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

ViT: Vision Transformer . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

ViT-OCR: ViT OCR . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

ViT-OCR-S: ViT OCR Small . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

ViT-OCR-B: ViT OCR Base . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

ViT-OCR-L: ViT OCR Large . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

ViT-OCR-H: ViT OCR Huge . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

# Chapter 1

# Introduction

This is the introduction chapter.