

Project

End-to-End Khmer OCR with Font Awareness

Abstract

Converting Khmer documents into editable formats like Markdown or HTML is challenged by the frequent use of multiple font styles within a single line—specifically Moul (for titles and honorifics) and Regular scripts. Current OCR solutions often lose this semantic formatting during extraction. We propose a unified Khmer OCR architecture that performs simultaneous character recognition and font classification using a shared CNN backbone. By branching a font classification head from the primary feature extractor, the model identifies styles (Regular, Moul, Bold, etc.) without additional computational overhead. This approach enables the generation of rich-text documents that preserve the visual hierarchy and semantic intent of the original Khmer source.

Introduction

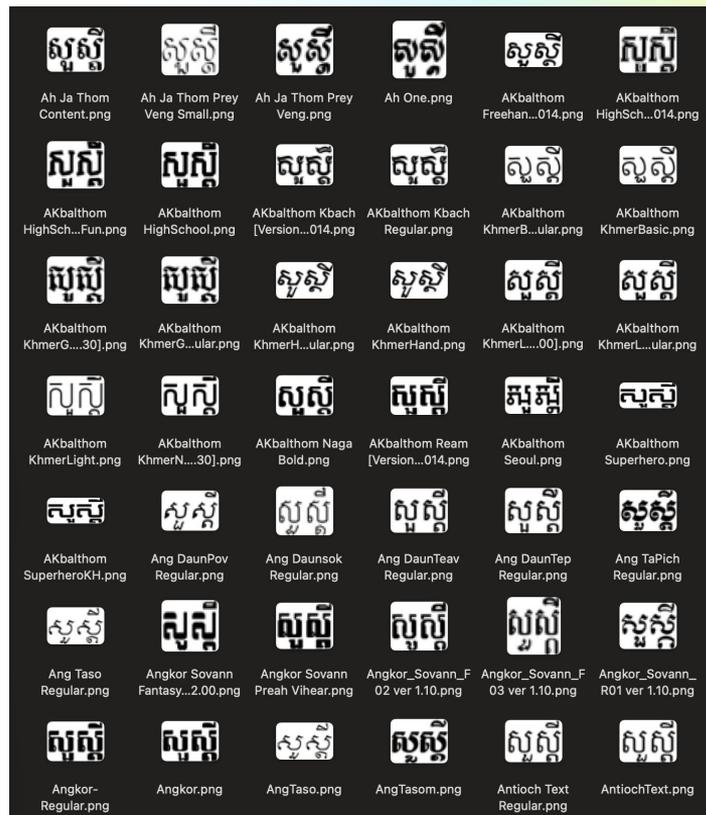
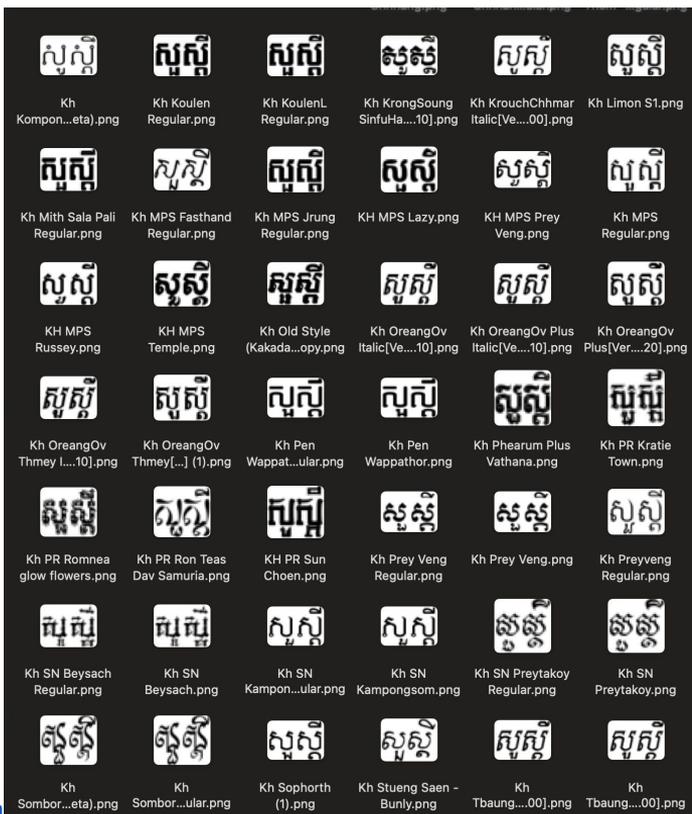
The digital transformation of Khmer printed documents into structured formats (PDF to Docx/Markdown) requires more than simple character accuracy. A defining characteristic of Khmer typography is the functional use of distinct font styles. For instance, **Moul** fonts are traditionally reserved for headings, names of high-ranking officials, and sacred texts, while Regular fonts are used for body content.

Standard OCR pipelines typically treat font variation as "noise" to be normalized, resulting in a loss of the document's original structure and emphasis. To bridge this gap, we introduce a multi-task learning architecture that utilizes a shared CNN backbone for both sequence recognition and font identification.

Dataset

To support the development of a robust Khmer Optical Character Recognition (OCR) system with Automatic Font Recognition (AFR), a large-scale synthetic text-line dataset was constructed. Due to the limited availability of publicly annotated Khmer OCR corpora, the dataset was generated using an automated rendering pipeline designed to capture both linguistic and typographic diversity inherent in the Khmer script. The total images are 2M+ sentences in Khmer language.

Dataset



Architecture & Configuration

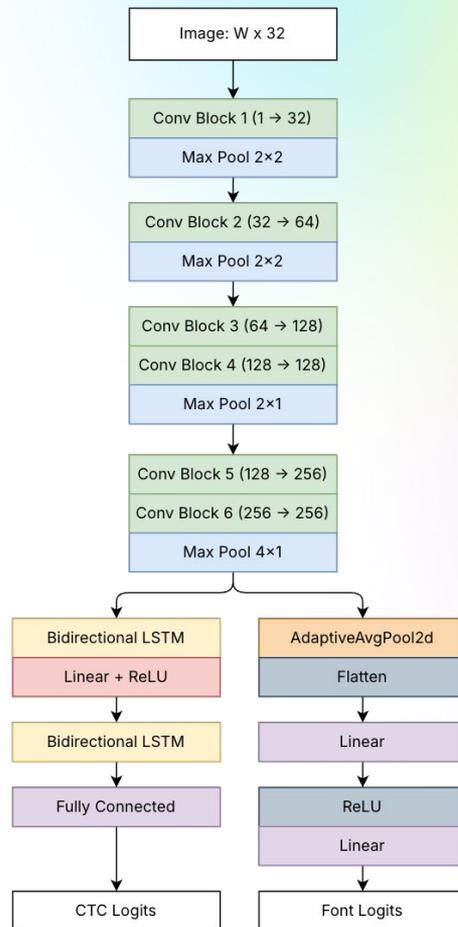
The backbone consists of the 6 convolution blocks followed two Bidirectional LSTMs for OCR task and AdaptiveAvgPool2d followed by two linear layers and ReLU for fonts classification head.

Optimizer: **AdamW**

Learning Rate: **0.0001**

OCR Loss Function: **CTCLoss**

Font Classify Loss Function: **Cross Entropy Loss**



Data Processing / Tokenizer

Category	Tokens
Special	[PAD], [SOS], [EOS], [UNK]
Consonants	ក, ខ, គ, ឃ, ង, ច, ឆ, ជ, ឈ, ញ, ដ, ប, ខ, ឈ, ណ, ត, ថ, ទ, ធ, ស, ប, ផ, ព, ភ, ម, យ, រ, ល, វ, ក, ម, ស, ហ, ឡ
Independent Vowels	អ, ម, អ, ត, ណ្ណ, ឧ, ឧ, ឌី, ឫ, ឫ, ឮ, ឮ, ង, ឮ, ឱ, ឱ, ឱ, ឌី

Dependent Vowels	ា, ិ, ី, ឺ, ឺ, ុ, ូ, ួ, ើ, ឿ, ៀ, េ, ៃ, ៃ, ោ, ៅ
Diacritics & Marks	ំ, ះ, ះ, ំ, ំ, ំ, ំ, ំ, ំ, ំ, ំ, ំ, ំ
Punctuation	។, ្រ, ្រ, ្រ, ្រ, ្រ, ្រ, ្រ
Numerals	០, ១, ២, ៣, ៤, ៥, ៦, ៧, ៨, ៩

OCR Evaluation

The model achieved **CER** of **0.0051** on the self-split validation set with the validation loss of **0.0167**.

Font Classification Evaluation

Font	Precision	Recall	F1-Score
<i>Regular</i>	0.8112	0.8690	0.8391
<i>Italic</i>	0.7691	0.8136	0.7908
<i>Bold</i>	0.8517	0.7933	0.8215
<i>BoldItalic</i>	0.7911	0.7623	0.7764
<i>Moul</i>	0.8182	0.9020	0.8580
<i>MoulLight</i>	0.9034	0.5325	0.6701

Result

ដឹកនាំធ្វើបទបង្ហាញជូនដល់សមាជិកអង្គទូតបរទេស និងអង្គការអន្តរជាតិនៅព្រះរាជាណាចក្រកម្ពុជា ពីការវិវត្តថ្មីៗ នៃ

Predicted Text: ដឹកនាំធ្វើបទបង្ហាញជូនដល់សមាជិកអង្គទូតបរទេសនិងអង្គការអន្តរជាតិនៅព្រះរាជាណាចក្រកម្ពុជាពីការវិវត្តថ្មីៗនៃ
Font: **Regular** (93%)

ស្តីពីការវិវត្តចុងក្រោយនៃស្ថានការណ៍ព្រំដែនកម្ពុជា-ថៃ

Predicted Text: ស្តីពីការវិវត្តចុងក្រោយនៃស្ថានការណ៍ព្រំដែនកម្ពុជាថៃ
Font: **Moul** (97%)

នាតិពាលមិនដឹងគួរ គ្មានគេសួរសោកចង់ដាក់ ឆ្លើយឆ្ងល់ដល់រាត់ទាត់ គួរហិនលក្ខណ៍ឆ្លាត់លើខ្លួន ។

Predicted Text: ជាតិពាលមិនដឹងគួរគ្មានគេសួរសោកចង់ដាក់ឆ្លើយឆ្ងល់ដល់រាត់ទាត់គួរហិនលក្ខណ៍ឆ្លាត់លើខ្លួន។
Font: **Italic** (69%)

ឈ្មោះក្តៅក្តួចក្រសែរ ដូចស្រាតកាយា បង្ហាញញាតិ ឈ្មោះក្តៅក្តួចសខ្ទមជាតិ ដូចលា តក់ណាប់បង្ហាញចោរ ។

Predicted Text: ឈ្មោះក្តៅក្តួចក្រសែរដូចស្រាតកាយាបង្ហាញញាតិឈ្មោះក្តៅក្តួចសខ្ទមជាតិដូចលាតក់ណាប់បង្ហាញចោរ។
Font: **MoulLight** (73%)

Text & Image Detection

We trained a YOLOv11 based text and image detection with the 5000+ Synthetic Documents generated by [Sone.js](https://github.com/seanghay/sones) (<https://github.com/seanghay/sones>)

១៤ ថ្ងៃ ២០២១ រៀង ២០២១

ពាក់ដីនគរក្រសួងសុខាភិបាលបានថ្លែងថា ក្រសួងសង្កេតឃើញមានការខូចខាតក្នុងការអនុវត្តនូវវិធានការសុខាភិបាល ក្នុងការទប់ស្កាត់ ការចម្លង និងការរីករាលដាលនៃជំងឺកូវីដ-១៩ ក្នុងអំឡុងពេលកាន់បិណ្ឌ និងបុណ្យភ្ជុំបិណ្ឌ។

លោកស្រីប្រជុំណតិ គី វណ្ណឌីន រដ្ឋលេខាធិការ និងជាអ្នកនាំពាក្យក្រសួងសុខាភិបាល បានថ្លែងថា «ដោយសារមានការសង្កេតឃើញមានការខូចខាតក្នុងការអនុវត្តវិធានការសុខាភិបាលក្នុងអំឡុងពេលកាន់បិណ្ឌ និងបុណ្យភ្ជុំបិណ្ឌនេះ ក្នុងនាមក្រសួងសុខាភិបាល ខ្ញុំសូមធ្វើការអំពាវនាវឱ្យប្រជាពលរដ្ឋទាំងអស់ ជាពិសេស លោកគារាមចារី លោកគារាមយាយ គណៈគ្រប់គ្រងសាលាពុទ្ធិក និងគណៈកម្មការគ្រប់គ្រងទីក្រុងភ្នំពេញ ទៅទាំងប្រទេស ព្រមទាំងពុទ្ធបរិស័ទទាំងអស់ ត្រូវបន្តអនុវត្តវិធានការអនាម័យការពារខ្លួន ពិការចម្លងវីរុសកូវីដ-១៩ ការចម្លងគ្រប់គ្រង និងការចម្លងជំងឺឆ្លងផ្សេងៗទៀត»។

លោកស្រីបានបន្តថាបង្កបង្កើនប្រសិទ្ធភាព អនុវត្តវិធានការសាមញ្ញៗមួយចំនួន ដូចជា រក្សាអនាម័យដៃស្អាតជានិច្ច ដោយលាងសម្អាតដៃក្រៃបន្តបន្ត និងកុំយកដៃដៃ មិនទាន់បានលាងសម្អាត ទៅប៉ះផ្ទៃមុខ មាត់ ឬចម្លោះ និងគ្រវែង ប្រើសារីទឹក កៅដៃ ក្រដាស និងចង្កា ប្រើប្រាស់ប្រាក់ចម្លង លោកស្រីបានបន្តថា «ចំណុចទី ២ ពាក់ម៉ាស់ ក្នុងករណីចាំបាច់ ជាពិសេសគ្រវែងពាក់ម៉ាស់ជាចាំបាច់ នៅពេលចូលទៅកាន់ទីក្រុងភ្នំពេញ ជួបជុំពុទ្ធបរិស័ទ ជាពិសេសមិនគ្រវែងជុំគ្នាទាំងប្រទេស (ការប្រគេនចង្កា និងទេយ្យគ្រឿងផ្សេងៗទៀត) ទៅតាមក្រុមគ្រួសារនីមួយៗ និងពេលសម្លេងព្រះងឹម»។

៣. រក្សាគម្លាតសុវត្ថិភាពបុគ្គលយ៉ាងតិចមួយបួនម៉ែត្រ ទៅបួនម៉ែត្រ មួយម៉ែត្រកន្លះឡើងទៅ ក្នុងនោះផងដែរ លោកស្រីបានបន្តថា «ជាពិសេសព្រះសង្ឃ ឬ បុគ្គលដែលឈរ មានរោគសញ្ញាផ្លាស្ស ឬមិនស្រួលខ្លួន ត្រូវសម្រាក និងនៅដាច់ដោយឡែក ដោយមិនត្រូវចូលរួមក្នុងពិធីជួបជុំក្នុងសាលានោះ ឬទទួលទេយ្យគ្រឿងក្នុងទីកន្លែងណាមួយនៃទីក្រុងភ្នំពេញ ដែលអាចបង្កឱ្យមានការចម្លងបន្តគ្នាពីមនុស្សម្នាក់ទៅមនុស្សម្នាក់»។



សម្តេចតេជោ អំពាវនាវប្រជាពលរដ្ឋបន្តចូលរួមទប់ស្កាត់ការឆ្លងរីករាលដាលជំងឺកូវីដ-១៩
១៦ ខែ កក្កដា ឆ្នាំ ២០២១ រៀង ០៣:៣៣

សម្តេចតេជោ ហ៊ុន សែន នាយករដ្ឋមន្ត្រីបានអំពាវនាវឱ្យប្រជាពលរដ្ឋបន្តចូលរួមទប់ស្កាត់ការឆ្លងរីករាលដាលជំងឺកូវីដ-១៩ ។

សម្តេចតេជោ បានសរសេរលើបណ្តាញសង្គម នៅថ្ងៃទី១៤ ខែកញ្ញាថា អនុរមានាពិធីបុណ្យភ្ជុំបិណ្ឌ «បិណ្ឌ១២» ។

សម្តេចបញ្ជាក់ថា « ថ្ងៃទី១៤ ខែកញ្ញានេះ គឺជាថ្ងៃកាន់បិណ្ឌ១២ នៃពិធីបុណ្យភ្ជុំបិណ្ឌ ដែលជាពិធីបុណ្យប្រពៃណីជាតិរបស់ខ្មែរយើង ដែលមានតាំងពីបុរាណមក ។

នេះជាពេលវេលាមួយ ដែលបងប្អូនប្រជាពលរដ្ឋពុទ្ធបរិស័ទមានឱកាសក្នុងការធ្វើបុណ្យទាន សន្យុកសល និងការជួបជុំសាច់ញាតិបងប្អូន នៅតាមទីកន្លែងធិត ត្រាយ ជាពិសេសតាមទីអារាមនានា។

ពិធីបុណ្យកាន់បិណ្ឌ ឆ្នាំ២០២១ នេះ គឺប្រព្រឹត្តទៅចាប់ពីថ្ងៃ១៣ ដល់ថ្ងៃ ១៤ រាច ខែកុម្ភៈ គឺចាប់ពីថ្ងៃទី១៤ ដល់ថ្ងៃទី១៥ ខែ កញ្ញា ឆ្នាំ២០២១។

ជាមួយគ្នានោះដែរ សម្តេចក៏សូមអញ្ជើញគ្រប់លំដាប់ថ្នាក់ បង្កើនកិច្ចការពារសុវត្ថិភាព សន្តិសុខ សណ្តាប់ធ្នាប់សាធារណៈ រៀបរៀងប្រយោជន៍ និងសម្រួលចរាចរណ៍ឱ្យបានល្អប្រសើរ ដើម្បីបង្កលក្ខណៈងាយស្រួលជូនដល់បងប្អូនប្រជាពលរដ្ឋនិរទេសនៅក្នុងកំឡុងពេលបុណ្យភ្ជុំបិណ្ឌនេះ ។

សម្តេចបានគូសបញ្ជាក់ថា សូមបងប្អូនប្រជាពលរដ្ឋបន្តចូលរួមទប់ស្កាត់ការឆ្លងរីករាលដាលជំងឺកូវីដ-១៩ ដោយស្តាប់តាមការណែនាំរបស់រដ្ឋាភិបាល ជាពិសេស ទឹកអាប៉េកុល និងពាក់ម៉ាស់ក្នុងពេលចូលរួមពិធីដែលមានអ្នកចូលរួមច្រើន ដូចជា ក្នុងពិធីបុណ្យភ្ជុំបិណ្ឌនេះផងដែរ ។

កុំធ្វើស្រប្រហែលសម្រាប់ជំងឺកូវីដ-១៩ អោយសោះ ព្រោះជាពិធីឆ្លងរាតត្បាតសហប្រយោជន៍ ហើយក៏ទាន់មានផ្ទៃក្រហមនៅឡើយទេ ។

សូមជូនពររបស់បងប្អូនជូនជោគជ័យក្នុងកំឡុង និងក្រៅប្រទេសជួបតែសេចក្តីសុខសេចក្តីចម្រើនគ្រប់គ្នា ។



Full OCR

Image

រាជធានីភ្នំពេញ ៖ ដើមឈើមានអាយុកាលជាច្រើនឆ្នាំមួយដើម អត់ស្គាល់ថា ជាដើមអ្វីនោះទេ ត្រូវបានជនមិនស្គាល់ អត្តសញ្ញាណលួចកាប់ចោល បង្កការភ្ញាក់ផ្អើលដល់អាជ្ញាធរចុះទៅពិនិត្យ និងឃាត់មនុស្ស ៣នាក់យកទៅសួរនាំ កាលពីយប់ថ្ងៃទី ៨ ខែកុម្ភៈ ឆ្នាំ ២០២៦ នៅចំណុចផ្លូវ ៥៧ កែងផ្លូវ ៣៣៤ ក្នុងសង្កាត់បឹងកេងកងទី ១ ខណ្ឌបឹង កេងកង រាជធានីភ្នំពេញ ។

តាមប្រភពបានអោយដឹងថា ដើមឈើមួយដើមនេះ មានអាយុកាលជាច្រើនឆ្នាំមកហើយ ជាសម្បត្តិរបស់រដ្ឋផងដែរ លុះ នៅយប់ថ្ងៃកើតហេតុ ជនមិនស្គាល់អត្តសញ្ញាណបានធ្វើសកម្មភាពកាប់ដើមឈើនោះ មិនដឹងថាប៉ុន្មានយកទៅណានោះ ទេ រហូតដល់មានអ្នករាយការណ៍ជូនអាជ្ញាធរមូលដ្ឋានចុះទៅពិនិត្យ និងឃាត់មនុស្ស៣នាក់យកទៅសួរនាំ ។

ជុំវិញករណីខាងលើនេះ ខាងមន្ត្រីពាក់ព័ន្ធ មិនទាន់បានអោយដឹងថា អ្នកកាប់ដើមឈើរបស់រដ្ឋទាំងនេះ យកទៅធ្វើអ្វី អោយពិតប្រាកដនៅឡើយទេ ។



Prediction

រាជធានីភ្នំពេញ ៖ ដើមឈើមានអាយុកាលជាច្រើនឆ្នាំមួយដើម អត់ស្គាល់ថា ជាដើមអ្វីនោះទេ ត្រូវបានជនមិនស្គាល់ អត្តសញ្ញាណលួចកាប់ចោល បង្កការភ្ញាក់ផ្អើលដល់អាជ្ញាធរចុះទៅពិនិត្យ និង ឃាត់មនុស្ស ៣នាក់យកទៅសួរនាំ កាលពីយប់ថ្ងៃទី ៨ ខែកុម្ភៈ ២០២៦ នៅចំណុចផ្លូវ ៥៧ កែងផ្លូវ ៣៣៤ ក្នុង សង្កាត់បឹងកេងកងទី ១ ខណ្ឌបឹង កេងកង រាជធានីភ្នំពេញ ។

តាមប្រភពបានអោយដឹងថា ដើមឈើមួយដើមនេះ មានអាយុកាលជាច្រើន ឆ្នាំមកហើយ ជាសម្បត្តិរបស់រដ្ឋផងដែរ លុះ នៅយប់ថ្ងៃកើតហេតុ ជនមិនស្គាល់អត្តសញ្ញាណបានធ្វើសកម្មភាពកាប់ ដើមឈើនោះ មិនដឹងថាប៉ុន្មានយកទៅណានោះ ទេ រហូតដល់មានអ្នករាយការណ៍ជូនអាជ្ញាធរមូលដ្ឋានចុះទៅពិនិត្យ និងឃាត់ មនុស្ស៣នាក់យកទៅសួរនាំ ។ ជុំវិញករណីខាងលើនេះ ខាងមន្ត្រីពាក់ព័ន្ធ មិនទាន់បានអោយដឹងថា អ្នកកាប់ ដើមឈើរបស់រដ្ឋទាំងនេះ យកទៅធ្វើអ្វី អោយពិតប្រាកដនៅឡើយទេ ។

