

# ການເພີ່ມປະສິດທິພາບການຮູ້ຈຳຕົວອັກສອນປ້າຍລົດລາວຈາກ ຮູບພາບໂດຍອັດຕະໂນມັດ ດ້ວຍວິທີ ການ jTessBoxEditor ແລະ Tesseract

ວິມິນທາ ຂຽວວົງພະຈັນ<sup>1\*</sup>, ຄອນປະເສີດ ສຸນາເຄນ<sup>1</sup>, ສະຫວາດ ໄຊປະດິດ<sup>1</sup>, ວິນັດ ເມກທະນາວັນ<sup>2</sup>, ຂັນທະນຸ ຫລວງ  
ໄຊຊະນະ<sup>1</sup>, ຊາລິສາ ຈັນທະວົງ<sup>1</sup>

ພາກວິຊາວິສະວະກຳຄອມພິວເຕີ ແລະ ເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ, ຄະນະວິສະວະກຳສາດ, ມະຫາວິທະຍາໄລແຫ່ງຊາດລາວ,  
ສປປ ລາວ

\* ຕິດຕໍ່ຜົວຜົນ:

ວິມິນທາ ຂຽວວົງພະຈັນ, ພາກວິຊາ  
ວິສະວະກຳຄອມພິວເຕີ ແລະ ເຕັກໂນ  
ໂລຊີຂໍ້ມູນຂ່າວສານ, ຄະນະ  
ວິສະວະກຳສາດ, ມະຫາວິທະຍາໄລ  
ແຫ່ງຊາດລາວ, ເບີໂທ: +856 20  
5559 2414,  
ອີເມວ:

[vimontha@fe-nuol.edu.la](mailto:vimontha@fe-nuol.edu.la)

<sup>1</sup> ພາກວິຊາວິສະວະກຳຄອມພິວເຕີ  
ແລະ ເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ,  
ຄະນະວິສະວະກຳສາດ,  
ມະຫາວິທະຍາໄລແຫ່ງຊາດລາວ,  
<sup>2</sup> ຄະນະວິສະວະກຳສາດ,  
ມະຫາວິທະຍາໄລສຸພານຸວົງ

ຂໍ້ມູນບົດຄວາມ:

ການສົ່ງບົດ: 05 ກັນຍາ 2023  
ການປັບປຸງ: 20 ພະຈິກ 2023  
ການຕອບຮັບ: 05 ທັນວາ 2023

## ບົດຄັດຫຍໍ້

ການຮູ້ຈຳຕົວອັກສອນ ດ້ວຍແສງ OCR ເປັນເຕັກໂນໂລຊີ ສຳຄັນທີ່ແປງ  
ຂໍ້ຄວາມທີ່ພິມ ຫຼື ຂຽນດ້ວຍລາຍມືໃຫ້ເປັນຂໍ້ຄວາມທີ່ເຄື່ອງສາມາດອ່ານໄດ້ ເຮັດ  
ໃຫ້ສາມາດປະມວນຜົນ ແລະ ວິເຄາະຂໍ້ມູນໄດ້ໂດຍອັດຕະໂນມັດແຕ່ເຖິງຢ່າງໃດກໍ່  
ຕາມ OCR ສຳລັບ ພາສາທີ່ມີຂໍ້ຄວາມຊັບຊ້ອນ ເຊັ່ນພາສາລາວແມ່ນສິ່ງທີ່ທ້າທາຍ  
ເນື່ອງຈາກວ່າລັກສະນະຕົວອັກສອນທີ່ຊັບຊ້ອນ. ດັ່ງນັ້ນ, ບົດຄົ້ນຄວ້າວິໄຈໄດ້ນຳສະເ  
ໜີ ການໃຊ້ເຄື່ອງມື jTessBoxEditor ສຳລັບການຝຶກ OCR ໃຫ້ຮູ້ຈຳຂໍ້ຄວາມ  
ພາສາລາວ. ເຄື່ອງມື jTessBoxEditor ແມ່ນສ່ວນເສີມຂອງ ເຄື່ອງມື Tesseract  
OCR ມີສ່ວນໜ້າຕ່າງທີ່ໃຊ້ງານງ່າຍ ສຳລັບການສ້າງ ແລະ ປັບແຕງຂໍ້ມູນການຝຶກ.  
ບົດຄົ້ນຄວ້າວິໄຈໄດ້ສະແດງຂະບວນການສ້າງແບບຈຳລອງ OCR ສຳລັບຂໍ້ຄວາມ  
ພາສາລາວ ລວມເຖິງການເກັບຂໍ້ມູນ, ການສ້າງການຝຶກຂໍ້ມູນ ແລະ ການວັດປະເມີນ  
ຜົນການສ້າງແບບຈຳລອງ. ຜົນການທົດເຫັນວ່າ ການສ້າງຊຸດຂໍ້ມູນ, ການສ້າງໄຟລ  
Tiff/Box ໂດຍການແຍກຕົວອັນສອນ, ແຍກຕົວພະຍັນຊະນະ ຫຼື ການປະສົມ  
ພະຍັນຊະນະ ແລະ ຕົວອັກສອນ ການຮູ້ຈຳຕົວອັກສອນພາສາລາວໃນການປ້າຍລົດ  
ລາວແມ່ນມີປະສິດ ທິພາບຄືກັນ ໂດຍມີການຝຶກແບບ train with existing box  
ແລະ ການຝຶກແບບ shape clustering. ໃນສ່ວນຄວາມໄວ ແລະ ຄວາມບໍ່ຊັບຊ້ອນ  
ເຫັນວ່າ ການສ້າງຊຸດຂໍ້ມູນແບບການແຍກຕົວອັກສອນ, ແຍກພະຍັນຊະນະ ແລະ ໃຊ້  
ການຝຶກແບບ train with existing box ແມ່ນມີຜົນດີກວ່າ.

**ຄຳສັບສຳຄັນ:** ການຮູ້ຈຳຕົວອັກສອນ, ຕົວອັກສອນລາວ, ຊຸດຂໍ້ມູນ, ການຝຶກ, ປ້າຍ  
ລົດລາວ, ເຄື່ອງມື jTessBoxEditor, ການປະມວນຜົນຮູບພາບ

## Enhancing Lao Vehicle License Plate Recognition through jTessBoxEditor and Tesseract Optical Character Recognition (OCR)

Vimontha Khieovongphachanh<sup>1\*</sup>, Khonepaserth Sounakhen<sup>1</sup>, Savath Saypadith<sup>1</sup>, Vinath Mekthanavanh<sup>2</sup>, Khanthanou Luangxaysana<sup>1</sup>, Xalisa Chanthavong<sup>1</sup>

Department of Computer Engineering and Information Technology, Faculty of Engineering, National University of Laos, Lao PDR

### \*Correspondence:

Vimontha  
Khieovongphachanh,  
Department of Computer  
Engineering and  
Information Technology,  
Faculty of Engineering,  
National University of  
Laos, PDR

Tel: +856 20 5559 2414,

Email:

[vimontha@fe-nuol.edu.la](mailto:vimontha@fe-nuol.edu.la)

<sup>1</sup>Department of Computer  
Engineering and  
Information Technology,  
Faculty of Engineering,  
National University of  
Laos, Lao PDR

<sup>2</sup>Faculty of Engineering,  
Souphanouvong University

### Article Info:

Submitted: Sep 05, 2023

Revised: Nov 20, 2023

Accepted: Dec 05, 2023

### Abstract

Optical Character Recognition (OCR) is a critical technology that converts printed or handwritten text into machine-readable text, enabling automated data processing and analysis. However, OCR for languages with complex text, such as Lao, presents unique challenges due to the intricate nature of the characters. This paper introduces the use of the jTessBoxEditor tool for training an OCR engine to recognize Lao text. The jTessBoxEditor tool, an extension of the Tesseract OCR engine, provides a user-friendly interface for creating and refining training data. The experimental results demonstrate the effectiveness of various techniques in recognizing Lao characters on Lao car plate. These techniques include data set creation, the generation of Tiff/Box files through character separation, and the isolation or combination of consonants and letters. Notably, the utilization of both existing box training and shape clustering training contributes to improved recognition performance. Furthermore, our findings highlight the importance of speed and simplicity in the OCR modeling process. Specifically, the creation of a dataset that involves character separation and consonant isolation, coupled with the use of existing box training, emerges as the most efficient and effective approach for Lao text recognition.

**Keyword:** Optical Character Recognition (OCR), Lao Character, Data Set, Data training, Lao Car Plate, jTessBoxEditor tool, Image Processing.

### 1. ພາກສະເໜີ

ເຕັກໂນໂລຊີ OCR (Optical Character Recognition) ແມ່ນການແຍກຂໍ້ຄວາມອອກຈາກຮູບພາບ ຫຼື ການແປງຂໍ້ຄວາມທີ່ຢູ່ໃນຮູບພາບມາເປັນຂໍ້ຄວາມດິຈິຕອນໂດຍການສາຍແສງ ເຊິ່ງມີການລິເລີ່ມຕັ້ງແຕ່ປີ ຄ.ສ 1920 ໂດຍນັກຝຶກຊາວເຢຍລະມັນທີ່ມີຊື່ວ່າ Emanuel Goldberg ໄດ້ປະດິດອຸປະກອນທີ່ເປັນຈຸດເລີ່ມຕົ້ນຂອງເຕັກໂນໂລຊີ OCR ເຊິ່ງເປັນຕົວອ່ານຕົວອັກສອນແລ້ວປຽນເປັນລະຫັດໂທລະເລກມາດຕະຖານ (Phommachanh & Phommachanh, 2019). ຫຼັງຈາກນັ້ນໃນປີ ຄ.ສ 1927 Goldberg ກໍໄດ້ພັດທະນາເຕັກໂນໂລຊີ OCR ໃຫ້ມີຄຸນນະພາບດີກວ່າເກົ່າກັບດ້ວຍສິ່ງປະດິດທີ່ມີຊື່ວ່າ “Statistical Machine” ເຄື່ອງມືທີ່ເຮັດໜ້າທີ່

ຊອກຫາບ່ອນເກັບມ້ຽນໄມໂຄຣຟິມ (Micro film) ເຊິ່ງເຮັດໜ້າທີ່ໃນການລະບຸຂໍ້ມູນດ້ວຍການຈີ້ດ້ວຍແສງ ໂດຍການເຮັດວຽກຂອງເຄື່ອງອ່ານສະຖິຕິໄດ້ເຮັດໃຫ້ Goldberg ປະສົບຜົນສໍາເລັດ ແລະ ໄດ້ຮັບສິດທິບັດສໍາລັບການປະດິດນີ້ກ່ອນທີ່ຈະຖືກຊື້ໂດຍ IBM. ຈົນມາເຖິງປີ ຄ.ສ 2000 ທີ່ OCR ອອນໄລນ໌ (Web OCR) ຖືກເປີດໃຊ້ຜ່ານຄລາວ (Cloud) ແລະ ແອັບພລິເຄຊັນເທິງມືຖື (Mobile Application) ທີ່ເພີ່ມຄວາມສະດວກໃນການນໍາໃຊ້ໂດຍຜ່ານໂທລະສັບມືຖືທີ່ມີຄວາມສະຫຼາດຫຼາຍຂຶ້ນ ເຊິ່ງເຮັດໜ້າທີ່ໂດຍການແຍກຂໍ້ຄວາມທີ່ບັນທຶກໄວ້ໂດຍໃຊ້ກ້ອງຖ່າຍຮູບຂອງອຸປະກອນ ຖ້າອຸປະກອນບໍ່ຮອງຮັບ OCR ກໍຍັງມີວິທີການອ່ານຕົວອັກສອນໂດຍໃຊ້ OCR API ທີ່ເຮັດວຽກໂດຍການແຍກຂໍ້ຄວາມອອກຈາກໄຟລ໌

ຮູບພາບ, ຈາກນັ້ນ ມັນຈະສົ່ງຄືນຂໍ້ຄວາມທີ່ແຍກອອກມາ ໄປຫາແອັບພລິເຄຊັນອຸປະກອນເພື່ອດຳເນີນການຕໍ່ໄປ.

ເຕັກໂນໂລຊີ OCR (Kaur & Singh, 2019; Kumar & Singh, 2022) ມີຄວາມສຳຄັນຫຼາຍ ແລະ ຍັງ ໄດ້ຖືກນຳໃຊ້ຢ່າງຫຼວງຫຼາຍໃນຊີວິດປະຈຳວັນຂອງພວກ ເຮົາຍົກຕົວຢ່າງ: ເຄື່ອງສະແກນເອກະສານ, ແອັບແປພາສາ, ເຄື່ອງມືປ່ຽນຂໍ້ຄວາມຈາກຮູບພາບເປັນຂໍ້ຄວາມຕົວອັກ ສອນ ແລະ ອື່ນໆ ເຊິ່ງໃນເຕັກໂນໂລຊີ OCR ຫຼື ການຮູ້ຈຳ ອັກສອນພາສາລາວມີຄວາມທ້າທາຍ ເນື່ອງຈາກພາສາລາວ ມີເອກະລັກ ແລະ ໂຄງສ້າງທີ່ຊັບຊ້ອນ ບໍ່ຄືກັບຕົວອັກສອນ ພາສາອັງກິດ ເຊິ່ງການຄົ້ນຄວ້າວິໄຈ ແລະ ການນຳໃຊ້ ແມ່ນຍັງມີຄວາມຖືກຕ້ອງຂອງຂໍ້ມູນຍັງຕໍ່າ ແລະ ຍັງບໍ່ຖືກ ໄວຍະກອນໃນພາສາລາວ ເຮັດໃຫ້ຫຼາຍພາກສ່ວນຍັງບໍ່ສາ ມາດໃຊ້ເຄື່ອງມືທີ່ອ່ານວຍຄວາມສະດວກໃນການສ້າງ ຫຼື ແປເອກະສານ. ການຄົ້ນຄວ້າທີ່ຜ່ານມາ ສຳລັບພາສາທີ່ຊັບ ຊ້ອນ ເຊັ່ນ: ພາສາອາລາບິກ ຫຼື ພາສາຈີນ ເຊິ່ງຂໍ້ມູນການ ຝຶກທີ່ມີປະສິດ ທິພາບສູງ ແລະ ເຄື່ອງມື jTessBoxEditor (Singh, & Singh, 2018; Sharma, & Gupta, 2019) ໄດ້ຮັບໃນການສ້າງແບບຈຳລອງການຈຳຕົວອັກສອນທີ່ ຊັບຊ້ອນປະກອບກັບການຮູ້ຈຳອັກສອນພາສາລາວຍັງບໍ່ທັນ ມີການເຜີຍ ແຜ່. ດັ່ງນັ້ນ, ບົດຄົ້ນຄວ້າວິໄຈນີ້ ຈຶ່ງມີຈຸດປະ ສົງນຳໃຊ້ການຝຶກການຮູ້ຈຳຕົວອັກສອນລາວ ດ້ວຍ jTessBoxEditor ແລະ Serack Tesseract Trainer (Ammar & Koubaa, & Boulila & Benjdira & Alhabashi, 2023), (Yap, W. J., & Tan, C. C., 2019), (Hlaing, Z. T., Zaw, Z., & Oo, K. Z., 2018) ເພື່ອສ້າງແບບຈຳລອງ ການຮູ້ຈຳພາສາລາວ ທີ່ ຖືກຕ້ອງສຳລັບຕົວອັກສອນ ໂດຍ ປະຍຸກໃຊ້ກັບ ບ້າຍລິດ ໃນປະເທດລາວ.

## 2. ອຸປະກອນ ແລະ ວິທີການ

### 2.1 ເຄື່ອງມື jTessBoxEditor

jTessBoxEditor ເປັນເຄື່ອງມືຕົວຊ່ວຍທີ່ມີໜ້າ ຕ່າງກຣາຟິກສຳລັບຜູ້ໃຊ້ (Graphical User Interface: GUI) ເຊິ່ງເປັນເຄື່ອງມືຊ່ວຍໃນການສ້າງ, ແກ້ໄຂ ແລະ ປັບ ປຸງຂໍ້ມູນການຝຶກສຳລັບ Tesseract OCR (Khan & Khan, 2019) ໂດຍ Tesseract ເປັນເຄື່ອງມື OCR ແບບ Open Source ອີກຕົວໜຶ່ງທີ່ມີຄວາມສາມາດໃນການຮູ້ຈຳ ພື້ນຖານ ໃນຂະນະທີ່ jTessBoxEditor ເປັນເຄື່ອງມືຕົວ ຊ່ວຍທີ່ມີໜ້າຕ່າງກຣາຟິກສຳລັບຜູ້ໃຊ້ (Graphical User Interface: GUI) ແບບ Open source ທີ່ໃຊ້ສຳລັບແກ້ ໄຂໄຟລ໌ຂໍ້ມູນການ ການຝຶກ ທີ່ໃຊ້ໃນລະບົບການຮູ້ຈຳ ໄດ້ ຮັບການອອກແບບມາໂດຍສະເພາະສຳລັບການເຮັດວຽກ

ຮ່ວມກັບ Tesseract OCR ເຊິ່ງເປັນເຄື່ອງມື OCR ທີ່ ພັດທະນາໂດຍ Google jTessBoxEditor ຊ່ວຍໃຫ້ຜູ້ໃຊ້ ສາມາດແກ້ໄຂໄຟລ໌ບ່ອກ (Box) ເຊິ່ງມີຂໍ້ມູນກ່ຽວກັບຕຳແ ໜ່ງ ແລະ ຂະໜາດຂອງລັກສະນະ ຫຼື ອັກສອນແຕ່ລະຕົວໃນ ຮູບພາບໄຟລ໌ບ່ອກເຫຼົ່ານີ້ຈຳເປັນສຳລັບການຝຶກ Tesseract OCR ໃຫ້ຮູ້ຈຳຂໍ້ຄວາມໃນຮູບພາບໄດ້ຢ່າງຊັດເຈນ ດ້ວຍ jTessBoxEditor ຜູ້ໃຊ້ສາມາດເປີດໄຟລ໌ບ່ອກ (Box) ສະແດງພາບທີ່ກ່ຽວຂ້ອງ ແລະ ທຳການປັບປ່ຽນ Box ດ້ວຍ ຕີນເອງຖ້າຈຳເປັນ. Tesseract ຖືກພັດທະນາໂດຍໂປຣ ແກຣມເມີ ທີ່ມີຊື່ວ່າ Serak Shiferaw ເຊິ່ງໂປຣແກຣມ ເປັນໂປຣແກຣມສຳລັບ ຝຶກ OCR Model (Phommachanh & Phommachanh, 2020) ແບບ ສຳເລັດຮູບຈາກໄຟລ໌ Tiff/Box ພ້ອມທັງກຳນົດອັກສອນ ສຳລັບ Train ໄດ້ນຳອີກ ເຊິ່ງການ Train Model ສຳເລັດ ຈະໄດ້ ໄຟລ໌ Model ສຳລັບໃຊ້ກັບ Tesseract OCR. ນຳ ໃຊ້ງານ Serak Tesseract trainer ໃນການນຳໄຟລ໌ Tiff/Box ມາ Train ເປັນ OCR Model ແລະ ຕັ້ງຄ່າຕົວ ອັກສອນຕ່າງໆ ຂອງ Model. ເຊິ່ງຖືກນຳໃຊ້ໃນຕົວສອນ ລາຕິນ.

### 2.2 ວິທີການ

ແມ່ນການສຸມບ້າຍລິດລາວ 100 ບ້າຍ ເພື່ອທຳການ ອ່ານໃຫ້ເປັນຕົວເລກ ແລະ ຕົວອັກສອນ ແລະ ທົດສອບ ຄວາມແມ້ນຍ້າຄວາມຖືກຕ້ອງຂອງອ່ານບ້າຍລິດລາວ.

ຂັ້ນຕອນການທົດລອງ ດັ່ງ ຮູບພາບທີ 2 ຂັ້ນຕອນ ການຈົດຈຳຕົວເລກ ແລະ ຕົວອັກສອນ ບ້າຍລິດລາວ.

#### 2.2.1 ກຽມຊຸດຂໍ້ມູນ

- ຕົວຜະຍັນຊະນະລາວປະກອບມີທັງໝົດ: 27 ໂຕ
- ຕົວສະຫຼະລາວປະກອບມີທັງໝົດ: 27 ໂຕ
- ຕົວເລກປະກອບມີທັງໝົດ: 10 ໂຕ

#### 2.2.2 ຊຸດຂໍ້ມູນ

ຊຸດຂໍ້ ມູນສຳລັບໃຫ້ Machine Learning (Kumar & Kumar, 2018; Kaur & Singh, 2019) ຮຽນຮູ້ ແມ່ນຂໍ້ ມູນທີ່ ຈະນຳເຂົ້າລະບົບເພື່ອເຮັດໃຫ້ Machine Learning ຮຽນຮູ້ໂດຍລັກສະນະຂອງຂໍ້ມູນຈະ ປະກອບດ້ວຍຄຸນລັກສະນະ (Features) ແລະ ຄຳຕອບ ຂອງການຮຽນຮູ້ຊຸດຂໍ້ມູນນັ້ນ (Label).

- Label ແມ່ນຂໍ້ ມູນທີ່ເຮົາຕ້ອງການຈະຄາດເດົາ (Predict) ໃນ Linear Regression Model.
- Features ແມ່ນຄຸນລັກສະນະທີ່ ມີລັກສະນະສະເພາະຂອງຂໍ້ ມູນທີ່ເຮົາຕ້ອງການຈະຮຽນຮູ້

ໂດຍເຮົາຈະໃຊ້ຂໍ້ມູນເຫຼົ່ານີ້ໃນການ ເຝິກ ແລະ ການຄາດ ເດົາ.

### 2.2.3 ຂັ້ນຕອນການຝຶກ

ນຳໃຊ້ jstessboxeditor ສ້າງໄຟລ໌ Tiff/Box ຈາກ ຊຸດຂໍ້ມູນ ເຊິ່ງໄຟລ໌ນີ້ຈະເປັນການຕັດຄຳ.

ຫຼັງຈາກໄດ້ໄຟລ໌ Tiff/Box ແລ້ວກຳນົດມາຝຶກ ໃນ ໂປຣແກຣມ Serak Trainer ເພື່ອສ້າງ Model OCR ທີ່ ສາມາດໃຊ້ໄດ້ກັບຫຼາຍໆ Tools

ວິທີການສ້າງໄຟລ໌ Tiff/Box ຈະແບ່ງເປັນ 2 ຮູບ ແບບ:

- ຮູບແບບທີ1(Tiff/Box1): ການສ້າງ Tiff/Box ໂດຍແຍກແຕ່ລະຕົວພະຍັນຊະນະ ຫຼື ແຕ່ລະຕົວສະຫຼະ
- ຮູບແບບທີ2(Tiff/Box2): ການສ້າງ Tiff/Box ໂດຍການປະສົມພະຍັນຊະນະ ແລະ ຕົວສະຫຼະການຝຶກ Tiff/Box ຈະມີແບບ 2 ແບບ:
  - ການຝຶກແບບທີ 1: train with existing box
  - ການຝຶກແບບທີ 2: shape clustering

### 2.2.4 ການທົດລອງ

ສຳລັບຂັ້ນຕອນການທົດລອງແມ່ນເຮົາຈະນຳຂໍ້ມູນ ຮູບປ້າຍລົດ (ຈຳລອງ) ເຊິ່ງປ້າຍລົດທີ່ເປັນຮູບພາບຈະມີຂໍ້ ມູນທີ່ເປັນ Text ໄຟລ໌ເປັນໄຟລ໌ .txt ເປັນຂໍ້ມູນປ້າຍທີ່ເປັນ ຕົວອັກສອນຕົ້ນສະບັບ ແລະ ຈະມີຂໍ້ມູນປ້າຍລົດທີ່ເປັນ ໄຟລ໌ PNG ເຊິ່ງເຮົາຈະນຳເອົາໄຟລ໌ ຂໍ້ມູນທີ່ເປັນໄຟລ໌ PNG ເຂົ້າຂະບວນການ OCR (Optical Character Recognition) ເປັນເຊິ່ງຈະນຳເອົາຂໍ້ມູນທີ່ເປັນຮູບພາບມາ ເຂົ້າຂະບວນການ OCR (Optical Character Recognition) ເປັນຂະບວນການທີ່ມີການຈັດຊຸດຂໍ້ມູນ ແບ່ງອອກເປັນ 2 ແບບ: Tiff/ Box1 (ຂໍ້ມູນທີ່ມີຮູບແບບ ຮູບ ແບບເປັນການແຍກສະຫະ, ພະຍັນຊະນະ ແລະ ຕົວ ເລກ) ແລະ Tiff/Box2 (ຂໍ້ມູນທີ່ມີການປະສົມພະຍັນ ຊະນະ ແລະ ສະຫະ) ແລ້ວນຳຝຶກ. ໂດຍການຝຶກຈະປະກອບ

### 2.2.6 ການທົດສອບຄວາມແມ້ນຢຳ

ການຄິດໄລ່ຄວາມແມ້ນຢຳ ແມ່ນໃຊ້ວິທີ Levenshtein Distance Equation (Nam, 2018)

$$lev_{a,b}(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) = 0 \\ \min\{lev_{a,b}(i - 1, j) + 1, lev_{a,b}(i, j - 1) + 1, lev_{a,b}(i - 1, j - 1) + 1_{(a_i \neq b_j)}\} & \text{otherwise} \end{cases} \quad (1)$$

$$accuracyrate = (1 - \frac{lev}{(\text{length}(a), \text{length}(b))}) * 100 \quad (2)$$

a: ຕົວອັກສອນຕົ້ນສະບັບ

j: ຕຳແໜ່ງຕົວອັກສອນຈາກຂະບວນການ OCR

b: ຕົວອັກສອນຈາກຂະບວນການ OCR

i: ຕຳແໜ່ງຕົວອັກສອນຕົ້ນສະບັບ

ມີ 2 ຮູບແບບຄື: Train with existing box ແລະ Shape Clustering ເຊິ່ງເມື່ອຜ່ານຂະບວນການ OCR (Optical Character Recognition) ເຮົາກໍຈະໄດ້ຮັບຂໍ້ມູນອັກສອນ ທີ່ຢູ່ໃນຮູບພາບໂດຍຈະເປັນອັກສອນດິຈິຕອນທີ່ຜ່ານ ຂະບວນການມາແລ້ວ ແລະ ນຳເອົາຕົວອັກສອນທີ່ຜ່ານການ ຝຶກທັງ 2 ຮູບແບບມາປຽບທຽບກັບອັກສອນຕົ້ນສະບັບ.

### 2.2.5 ຜົນການທົດລອງ

ຜົນການທົດລອງຂອງການຝຶກຮູບແບບ Train with existing box ແລະ Shape Clustering ແລະ ນຳ ໃຊ້ຂໍ້ມູນທັງສອງຊຸດເຊິ່ງຂໍ້ມູນທີ່ນຳມາຝຶກປະກອບມີ: ພະຍັນຊະນະ, ສະຫຼະ ແລະ ຕົວເລກ ຄື Tiff/Box1 (ຂໍ້ມູນ ທີ່ມີຮູບແບບເປັນການແຍກສະຫຼະ, ພະຍັນຊະນະ ແລະ ຕົວ ເລກ) ແລະ Tiff/Box2 (ຂໍ້ມູນທີ່ມີການປະສົມ ພະຍັນຊະນະ ແລະ ສະຫຼະ) ຜ່ານຂະບວນການ OCR ອີງ ຕາມຕາຕະລາງທີ 4 ສັງເກດໄດ້ຜ່ານຕາຕະລາງຈະປະກອບ ມີ ຕົວອັກສອນຕົ້ນແບບ 20 ຊຸດທີ່ເຮົານຳມາໃຊ້ໃນການ ຝຶກ ແລະ ຜົນລັບຂອງການຝຶກທັງສອງແບບໂດຍນຳໃຊ້ຊຸດ ຂໍ້ມູນທັງສອງແບບ Tiff/Box1 ແລະ Tiff/Box2 ລວມມີ ຜົນລັບທັງໝົດ 4 ແບບ: Train with existing box (Tiff/Box1), Train with existing box (Tiff/Box2), Shape Clustering (Tiff/Box1), Shape Clustering(Tiff/Box2) ເຊິ່ງແຕ່ລະແບບຈະໃຫ້ຜົນຮັບທີ່ ແຕກຕ່າງກັນ. ສຳລັບຂໍ້ຄວາມທີ່ຜ່ານການຝຶກມາແລ້ວແຕ່ ຜົນຮັບທີ່ຍັງບໍ່ເປັນທີ່ພິພິຈາລະນາມີ 5 ຊຸດຂໍ້ມູນ: ຫຼວງນ້ຳ ທາ (2), ຜັງສາລີ (2), ສາລະວັນ ເນື່ອງຈາກວ່າຂໍ້ຄວາມດັ່ງ ກ່າວມີຄວາມຊັບຊ້ອນຈຶ່ງເຮັດຜົນລັບທີ່ອອກມາບໍ່ດີປານ ໃດ. ສ່ວນຊຸດຂໍ້ມູນອື່ນໆແມ່ນກໍໄດ້ຜົນລັບທີ່ໃກ້ຄຽງກັບ ຊຸດຂໍ້ມູນຕົ້ນແບບ ແຕ່ກໍຍັງມີບາງຕົວອັກສອນ, ສະຫຼະ, ພະຍັນຊະນະ ທີ່ຍັງສະກົດຜິດຢູ່ບາງສ່ວນ. ສຳລັບຕົວເລກ ແມ່ນໄດ້ຮັບຜົນລັບທີ່ດີຫຼາຍຊຸດຂໍ້ມູນແຕ່ລະຊຸດທີ່ຜ່ານການ ຝຶກແລ້ວມີຄວາມແນ່ນອນຫຼາຍທີ່ສຸດລ້ວນແຕ່ແມ່ນ ຂໍ້ຄວາມທີ່ເປັນຕົວເລກ.

- Levenshtein Distance (lev): ຈຳນວນຄວາມແຕກຕ່າງລະຫວ່າງ ຕົວອັກສອນຕົ້ນສະບັບ ແລະ ຕົວອັກສອນຈາກຂະບວນການ OCR.

- accuracyrate: ຄວາມແມ້ນຍຳລະຫວ່າງ ຕົວອັກສອນຕົ້ນສະບັບ ແລະ ຕົວອັກສອນຈາກຂະບວນການ OCR ຄິດໄລ່ ເປັນເປີເຊັນ.

### 3. ຜົນໄດ້ຮັບ

ຜົນຂອງການຝຶກແບບທີ 1 ຈະເຫັນໄດ້ວ່າ Tiff/Box ຮູບແບບທີ 1 ຈະມີຄ່າຄວາມແມ້ນຍຳສູງກວ່າ Tiff/Box ຮູບແບບທີ 2 ທຸກປ້າຍລົດທີ່ນຳມາທົດລອງ ແລະ ມີຜົນຄ່າສະເລ່ຍຄວາມແມ້ນຍຳເປັນເປີເຊັນ 60.32% ໃນຂະນະທີ່ການຝຶກ Tiff/Box ຮູບແບບທີ 2 ເຫັນວ່າ ຄ່າຄວາມແມ້ນຍຳຂອງ Tiff/Box 1 ແລະ Tiff/Box 2 ມີຄ່າໃກ້ຄຽງກັນທຸກປ້າຍລົດ ແລະ ມີຜົນຄ່າສະເລ່ຍຄວາມແມ້ນຍຳເປັນເປີເຊັນ 59.11% ເຊິ່ງເຫັນໄດ້ວ່າຄວາມແມ້ນຍຳຂອງຮູບແບບທີ 1 ຈະສູງກວ່າ ແລະ ບໍ່ມີຂັ້ນຕອນໃນການຝຶກຂໍ້ມູນທີ່ຊັບຊ້ອນ ຖ້າການທົດລອງມີການຝຶກຂໍ້ມູນຫຼາຍຂຶ້ນກໍ່ຈະເຮັດໃຫ້ຜົນຂອງຮູບແບບທີ 1 ມີຄວາມແມ້ນຍຳສູງຂຶ້ນ.

- ການຝຶກແບບທີ 1: train with existing box ໃນTiff/Box 2 ຮູບແບບ. ດັ່ງຕາຕະລາງທີ 2

- ການຝຶກແບບທີ 2: Shape Clustering ໃນ Tiff/Box 2 ຮູບແບບ. ດັ່ງຕາຕະລາງທີ 3

### 4. ວິພາກຜົນ

ການຄົ້ນຄວ້າກ່ຽວກັບ OCR ຂອງພາສາລາວກໍ່ຄື ການຈື່ຈຳຕົວອັກສອນໃນຮູບແບບຂໍ້ຄວາມແມ້ນຍຳຍັງບໍ່ແຜ່ຫຼາຍ ແລະ ເຕັກນິກການຝຶກ Machine Learning ມີຮູບແບບເຕັກນິກທີ່ແຕກຕ່າງກັນເຊິ່ງການນຳໃຊ້ຂຶ້ນກັບຈຸດປະສົງທີ່ແຕກຕ່າງກັນ ແລະ ເຫັນວ່າແມ້ນຍຳຄວາມແມ້ນຍຳ 90% ຂຶ້ນໄປ ເຊິ່ງເຫັນວ່າມີຄວາມແມ້ນຍຳສູງ ຍ້ອນຂໍ້ມູນມີຂະໜາດໃຫຍ່ໃນການຝຶກ ແລະ ເຕັກນິກມີຄວາມຊັບຊ້ອນ ເຮັດໃຫ້ມີການປະມວນຜົນໃຊ້ເວລາຫຼາຍຂຶ້ນ (Kesom et al, 2018) (Beyene et al, 2019). ບົດຄົ້ນຄວ້ານີ້ສຶກສາເຕັກນິກທີ່ບໍ່ມີຄວາມຊັບຊ້ອນ ແລະ ການຝຶກຂໍ້ມູນບໍ່ຫຼາຍສຳລັບພາສາລາວ ເຊິ່ງເຕັກນິກທີ່ນຳມາໃຊ້ນີ້ແມ້ນຖືກນຳໃຊ້ກັບພາສາຕ່າງປະເທດ ເຊັ່ນ: ພາສາລາຕິນ ແລະ ພາສາພື້ນເມືອງອື່ນໆ ເຫັນວ່າ Tesseract OCR ມີຄວາມແມ້ນຍຳໃນລະດັບທີ່ໜ້າພໍໃຈ, ເຊິ່ງມີຄວາມແຕກຕ່າງກັບບົດຄົ້ນຄວ້າການນຳໃຊ້ເຄື່ອງມືໃຊ້ໃນການອ່ານເອກະສານພາສາລາຕິນ (Khadijaelgajoui et al., 2015)

ແມ້ນສາມາດປັບປຸງການອ່ານໄດ້ດີ ແຕ່ຍັງບໍ່ທັນໄດ້ລົງເລິກ ແລະ ຍັງບໍ່ທັນເພີ່ມ ລັກສະນະຂອງຟອນພາສາ. ບົດຄົ້ນຄວ້າອ່ານຂໍ້ຄວາມຈາກເອກະສານ PDF ໃນພາສາທີ່ມີຊັບພະຍາກອນຕ່ຳໂດຍໃຊ້ຕົວອັກສອນພື້ນເມືອງ, ບົດຄົ້ນ ຄວ້າໄດ້ປັບປຸງການປະຕິບັດຂອງ Tesseract, ເຄື່ອງມືການຮັບຮູ້ຕົວອັກສອນໂດຍການຝຶກມັນຢູ່ໃນຫຼາຍກວ່າ 20 ຕົວອັກສອນພື້ນເມືອງໃນ Tamil ແລະ Sinhala (Charangan Vasantharajan et al., 2019). ວິທີການສຳລັບການສະກັດຂໍ້ຄວາມຈາກເອກະສານ PDF ໂດຍໃຊ້ຕົວອັກສອນພື້ນເມືອງ. ຜູ້ຂຽນໄດ້ປັບປຸງການປະຕິບັດຂອງ Tesseract, ເຄື່ອງມືການຮັບຮູ້ຕົວອັກສອນ optical (OCR), ໂດຍການຝຶກອີບຣິມມັນຢູ່ໃນຫຼາຍກວ່າ 20 ຕົວອັກສອນພື້ນເມືອງໃນ Tamil ແລະ Sinhala. ວິທີການນີ້ຊ່ວຍຫຼຸດອັດຕາຄວາມຜິດພາດໃນລະດັບຕົວອັກສອນຂອງ Tesseract ແລະ ສະແດງໃຫ້ເຫັນວ່າ ການໃຊ້ຫຼັກການສຳລັບອັກສອນພື້ນເມືອງໃນ Tamil ແລະ Sinhala ແມ້ນສາທາດເຂົ້າໃຈຂໍ້ຄວາມ ແລະ ປັບປຸງຂະບວນການ OCR ແລະ ມີວິທີການຝຶກແບບໃໝ່. ເຄື່ອງອ່ານຂໍ້ຄວາມທີ່ຂຽນດ້ວຍມືຊ່ວຍຄົນພິການທາງສາຍຕາໂດຍການປ່ຽນຂໍ້ຄວາມທີ່ສະແກນ (hithra Selvaraj & Bhalaji Natarajan, 2018), ພິມ ແລະ ຂຽນດ້ວຍມືເປັນສຽງ. ເຄື່ອງສະແກນຫນ້າມືຖືກນຳໃຊ້ເພື່ອສະແກນຂໍ້ຄວາມ. ແອັບພລິເຄຊັນນີ້ໃຊ້ເຄື່ອງມື Tesseract OCR ເພື່ອອ່ານຂໍ້ຄວາມຈາກຮູບພາບ ແລະ ປ່ຽນເປັນສຽງເວົ້າ. ເຄື່ອງມື OCR ໄດ້ຖືກຝຶກໃຫ້ຮັບຮູ້ຂໍ້ຄວາມທີ່ຂຽນດ້ວຍມືສຳລັບ ພາສາອິນດູ ແລະ Bengali.

ຜ່ານການຝຶກທັງສອງຮູບແບບ ທັງ Train with existing box ແລະ Shape Clustering ສັງເກດໄດ້ວ່າຜົນຮັບທີ່ໄດ້ຮັບຄວາມແມ້ນຍຳດີທີ່ສຸດຈາກຮູບແບບ Train with existing box ແມ້ນການນຳໃຊ້ຊຸດຂໍ້ມູນ Tiff/Box1 (ຂໍ້ມູນທີ່ມີຮູບແບບເປັນການແຍກສະຫຼະ, ພະຍັນຊະນະ ແລະ ຕົວເລກ), ສ່ວນຜົນລັບທີ່ໄດ້ຮັບຄວາມແມ້ນຍຳດີທີ່ສຸດຂອງຮູບແບບ Shape Clustering ແມ້ນການນຳໃຊ້ຊຸດຂໍ້ມູນ Tiff/Box2 (ຂໍ້ມູນທີ່ມີການປະສົມພະຍັນຊະນະ ແລະ ສະຫຼະ). ຈາກຜົນລັບການປຽບທຽມຄວາມແມ້ນຍຳ ຈາກຮູບພາບທີ 10 ສັງເກດໄດ້ວ່າ ຮູບແບບ Train with existing box ທີ່ນຳໃຊ້ຊຸດຂໍ້ມູນ Tiff/Box1 (ຂໍ້ມູນທີ່ມີຮູບແບບເປັນການແຍກສະຫຼະ, ພະຍັນຊະນະ ແລະ ຕົວເລກ) ແມ້ນມີຄວາມແມ້ນຍຳຫຼາຍກວ່າ ຮູບແບບ Shape Clustering ທີ່ນຳໃຊ້ຊຸດຂໍ້ມູນ Tiff/Box2 (ຂໍ້ມູນທີ່ມີການປະສົມພະຍັນຊະນະ ແລະ

ສະຫຼະ) ບໍ່ວ່າຈະເປັນເປີເຊັນຄວາມແມ້ນຍ່າ ຫຼື ຜົນຮັບຜ່ານ ການຝຶກແມ້ນ ຮູບແບບ Train with existing box ທີ່ນຳ ໃຊ້ຊຸດຂໍ້ມູນ Tiff/Box1 ເຮັດໄດ້ດີກວ່າຫຼາຍ ເຖິງຈະໄດ້ ຮັບຊຸດຂໍ້ມູນທີ່ຊັບຊ້ອນແຕ່ກໍໄດ້ຜົນຮັບທີ່ດີ ແລະ ໃກ້ຄຽງ ກັບຊຸດຂໍ້ມູນທີ່ນຳມາຝຶກ. ຕົວຢ່າງ: ຕົວອັກສອນຕົ້ນສະບັບ (ເຊກອງ, ອັດຕະປື, ຫົວຜັນ, ບໍ່ແກ້ວ) ເມື່ອນຳອັກສອນຕົ້ນ ແບບມາທົດລອງໃນການຝຶກຮູບແບບ Train with existing box ໂດຍການນຳໃຊ້ຊຸດຂໍ້ມູນ Tiff/Box1 ຈະ ໄດ້ຜົນລັບດັ່ງນີ້ (ເຊກອງ, ອັດຕະປື, ຫົວຜັນ, ບໍ່ແກ້ວ), ສ່ວນຮູບແບບ Shape Clustering ທີ່ນຳໃຊ້ຊຸດຂໍ້ ມູນ Tiff/Box2 ແມ່ນໃຫ້ຜົນຮັບແມ້ນ (ເຊກອດີ, ອັດຕະປື, ຫົວ ຜັນ, ບໍ່ແກ້ວ). ຈຸດປະສົງຂອງບົດຄົ້ນຄວ້າແມ່ນເພື່ອນຳໃຊ້ ເຄື່ອງມືທີ່ບໍ່ມີຄວາມຊັບຊ້ອນໃນການຝຶກຂໍ້ມູນ ແລະ ການ ກະກຽມຂໍ້ມູນ ແລະ ສະພາບແວດລ້ອມເຫຼົ່ານີ້ແມ່ນຍັງບໍ່ມີ ງານວິໄຈໃນການຈື່ຈຳພາສາລາວ ສັງເກດເຫັນວ່າມີຄວາມ ແມ້ນຍ່າສູງເຖິງ 60.32% ແມ່ນເຮັດໄດ້ດີກວ່າຜົນໄດ້ຮັບ ຂອງຂໍ້ຄວາມແມ້ນຢູ່ໃນລະດັບທີ່ສາມາດອ່ານ, ເຂົ້າໃຈໄດ້ ແລະ ມີຄວາມໃກ້ຄຽງກັບຂໍ້ຄວາມຕົ້ນແບບທີ່ນຳມາຝຶກ.

### 5. ສະຫຼຸບ

ເຕັກໂນໂລຊີ Optical Character Recognition ຫຼື ການຮູ້ຈຳອັກສອນແມ່ນມີຄວາມທ້າທາຍ ເນື່ອງຈາກ ພາສາລາວມີເອກະລັກ ແລະ ໂຄງສ້າງທີ່ຊັບຊ້ອນ ບໍ່ຄືກັບ ຕົວອັກສອນ ພາສາອັງກິດ ເຊິ່ງການຄົ້ນຄວ້າວິໄຈ ໄດ້ນຳ ສະເໜີ ການຮູ້ຈຳຕົວອັກສອນ ພາສາລາວ ໂດຍນຳໃຊ້ ເຄື່ອງມືການຝຶກດ້ວຍ jTessBoxEditor ແລະ Serack Tesseract Trainer ແລະ ທົດລອງກັບການອ່ານປ້າຍລົດ ໃນປະເທດລາວທີ່ມີຕົວອັກສອນພາສາລາວ ແລະ ຕົວເລກ ໂດຍການສ້າງໄຟລ Tiff/Box, ການຕັດຄຳ 2 ຮູບແບບ ແລະ ແຕ່ລະຮູບມີການຝຶກແບບ train with existing box ແລະ shape clustering ຈາກຜົນການທົດລອງເຫັນ ວ່າ ການຕັດຄຳ ໂດຍແຍກແຕ່ລະຕົວພະຍັນຊະນະ-ແຕ່ລະ ຕົວສະຫຼະ ຫຼື ໂດຍການປະສົມພະຍັນຊະນະ ແລະ ຕົວສະຫຼະ ແມ່ນມີຜົນໄດ້ຮັບທີ່ດີແຕ່ຕ່າງການຝຶກ ສະຫຼຸບໄດ້ວ່າ ສາມາດນຳໃຊ້ Tiff/Box ແບບທີ່ແຍກໂຕພະຍັນຊະນະຈະ ມີຄວາມໄວ, ງ່າຍໃນການຝຶກ ແລະ ຈັດການຊຸດຂໍ້ມູນໄດ້ ດີກວ່າ. ໃນຕໍ່ໜ້າຈະມີການເພີ່ມຊຸດຂໍ້ມູນໃຫ້ຫຼາຍ ແລະ ອ່ານກັບປ້າຍລົດທີ່ຖືກຖ່າຍຈາກຕົວຈິງ ເພື່ອໃຫ້ມີຄວາມ ແມ້ນຍ່າ ແລະ ສາມາດນຳໃຊ້ໄດ້ທົ່ວໄປ.

### 6. ຂໍ້ຂັດແຍ່ງ

ຂ້າພະເຈົ້າໃນນາມຜູ້ຄົ້ນຄວ້າວິທະຍາສາດ ຂໍປະຕິ ຍານຕົນວ່າ ຂໍ້ມູນທັງໝົດທີ່ມີໃນບົດຄວາມວິຊາການດັ່ງ

ກ່າວນີ້ ແມ່ນບໍ່ມີຂໍ້ຂັດແຍ່ງທາງຜົນປະໂຫຍດກັບພາກສ່ວນ ໃດ ແລະ ບໍ່ໄດ້ເອື້ອປະໂຫຍດໃຫ້ກັບພາກສ່ວນໃດພາກ ສ່ວນໜຶ່ງ, ກໍລະນີມີການລະເມີດໃນຮູບການໃດໜຶ່ງ ຂ້າພະ ເຈົ້າມີຄວາມຍິນດີທີ່ຈະຮັບຜິດຊອບແຕ່ພຽງຜູ້ດຽວ.

### 7. ເອກະສານອ້າງອີງ

- Phommachanh, S., & Phommachanh, K. (2019). Development of Lao Vehicle License Plate Recognition System Using Tesseract OCR. *International Journal of Computer Applications*, 178(40), 1-5.
- Nam, E. (2018). Understanding the Levenshtein distance equation for beginners. *Medium*.
- Phommachanh, S., & Phommachanh, K. (2020). A Study on Lao Vehicle License Plate Recognition System Using Deep Learning. *International Journal of Innovative Technology and Exploring Engineering*, 9(3), 2228-2232.
- Singh, A., & Singh, S. (2018). Vehicle License Plate Recognition Using Deep Learning and OpenCV. *International Journal of Engineering Research & Technology*, 7(11), 337-341.
- Khan, M., & Khan, M. (2019). A Survey on License Plate Recognition Systems. *International Journal of Advanced Computer Science and Applications*, 10(2), 498-505.
- Ammar, A., Koubaa, A., Boulila, W., Benjdira, B., & Alhabashi, Y. (2023). A multi-stage deep-learning-based vehicle and license plate recognition system with real-time edge inference. *Sensors*, 23(4), 2120.
- Kumar, S., & Singh, S. (2022). Comparative analysis of EasyOCR and TesseractOCR for automatic license plate recognition using deep learning algorithm. In *2022 6th International Conference on Electronics, Communication and Aerospace Technology* (pp. 1-5). IEEE.
- Kaur, J., & Singh, S. (2019). License plate recognition using OpenCV and Tesseract OCR on Indian number plates. *International Journal of Engineering Research & Technology*, 8(12), 1-5.

- Li, X., Wang, Z., Wang, Y., & Zhang, H. (2018). Vehicle License Plate Recognition Based on Tesseract OCR Engine. In 2018 37th Chinese Control Conference (CCC) (pp. 8537-8541). IEEE.
- Vongpradhip, S., Sudsang, A., & Sinthupinyo, S. (2017). Enhanced Vehicle License Plate Recognition System Using Deep Learning. In 2017 14th International Joint Conference on Computer Science and Software Engineering (JCSSE) (pp. 1-5). IEEE.  
<https://doi.org/10.1109/JCSSE.2017.8026139>.
- Hlaing, Z. T., Zaw, Z., & Oo, K. Z. (2018). Vehicle License Plate Recognition System Based on Tesseract OCR Engine. In 2018 IEEE Conference on Computer Applications (ICCA) (pp. 1-4). IEEE.  
<https://doi.org/10.1109/COMAPP.2018.8377522>.
- Li, B., & Wang, C. (2019). Research on License Plate Recognition Based on Tesseract OCR. In 2019 5th International Conference on Control, Automation and Robotics (ICCAR) (pp. 320-323). IEEE.  
<https://doi.org/10.1109/ICCAR.2019.8813754>.
- Yap, W. J., & Tan, C. C. (2019). Vehicle License Plate Recognition System Based on Deep Learning and Tesseract OCR. In 2019 2nd International Conference on Computational Intelligence and Intelligent Systems (pp. 22-26). IEEE.  
<https://doi.org/10.1109/CIIS.2019.8883942>.
- Li, M., & Xu, Z. (2019). Vehicle License Plate Recognition System Based on Tesseract OCR Engine. In 2019 2nd International Conference on Frontiers of Control and Sensor Networks (CFCSN) (pp. 67-70). IEEE.  
<https://doi.org/10.1109/CFCSN47491.2019.00021>.
- Khadija El Gajoui, Fadoua Ataa Allah, Mohammed Oumsis. (2015). Training TESSERACT Tool for Amizig OCR. RECENT RESEARCHES in APPLIED COMPUTER SCIENCE: Proceedings of the 15th International Conference on Applied Computer Science (ACS15), Konya, Turkey May 20-22, 2015 (pp.172-179).
- Charangan Vasantharajan, Laksika Tharmalingam, Uthayasanker Thayasivam. (2019). Adapting the Tesseract Open-Source OCR Engine for Tamil and Sinhala Legacy Fonts and Creating a Parallel Corpus for Tamil-Sinhala-English.
- Chithra Selvaraj, Bhalaji Natarajan. (2018). Enhanced portable text to speech converter for visually impaired. International Journal of Intelligent Systems Technologies and Applications 17(1/2): 42.
- Kesom, K., & Phawapoothayanchai, P. (2018). Optical Character Recognition (OCR) enhancement using an approximate string matching technique. Engineering and Applied Science Research, 45(4), 282-289.
- Beyene, E. G. (2019). Handwritten and Machine printed OCR for Geez Numbers Using Artificial Neural Networks arXiv preprint arXiv:1911.06845.
- Barkov, A. (2019). The hundred-page machine learning book (Vol. 1, p. 32). Quebec City, QC, Canada: Andriy Burkov.

ຕາຕະລາງທີ 1. ຜົນການທົດລອງ

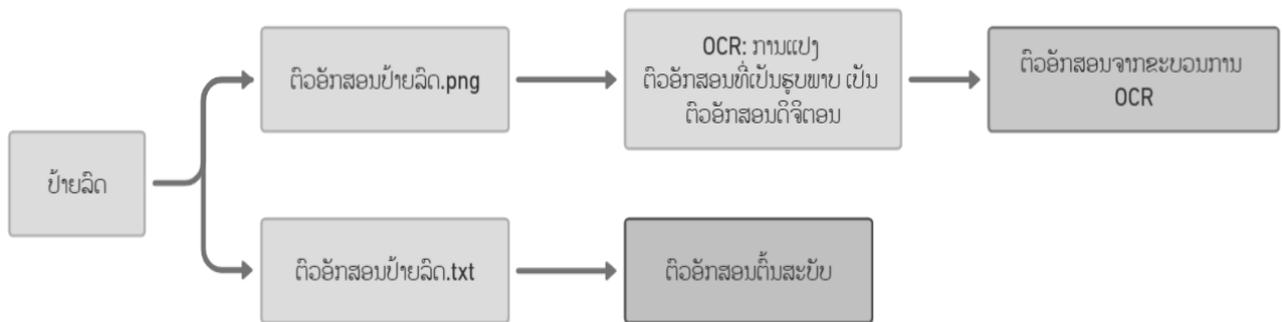
ລຳດັບ [ຮູບປ້າຍລົດ]	ຕົວອັກສອນຕົ້ນ ສະບັບ	ຕົວອັກສອນ ຈາກຂະບວນການ OCR			
		ການຝຶກແບບທີ1 Tiff/Box1	ການຝຶກແບບທີ1 Tiff/Box2	ການຝຶກແບບທີ2 Tiff/Box1	ການຝຶກແບບທີ2 Tiff/Box2
1	ເຊກອງ ຍປ 6155	ເຊກອງ ຍປ 6155	ດຸຊກອງ ຍປ 6155	ເຊກອດີ ຍັປ 6155	ເຊກອດີ ຍັປ 6155
2	ເຊກອງ ຂຜ 5514	ເຊກອງ ຂຜ 5514	ທຸຊກອງ ຂຜ 5514	ເຊກຍີ ຂຜ 5514	ເຊກຍີ ຂຜ 5514
3	ຫົວຜັນ ຖຜ 8732	ຫົວຜັນ ຖຜ 8732	ຫົວຜັນ ຖຜ 8732	ຫົວຜັນ ຖຜ 8732	ຫົວຜັນ ຖຜ 8732
.	.	.	.	.	.
.	.	.	.	.	.
.	.	.	.	.	.
19	ສາລະວັນ ຂຂ 9345	ສຸຖອອັນ ຂຂ 9345	ສຸຖອຽນ ຂຂ 9345	ສົງວຣິອັນ ຂຂ 9345	ສົງວຣິອັນ ຂຂ 9345
20	ອັດຕະປື ຍບ 5670	ອັດຕຸຣປື ຍບ 5670	ອັດຕຸຣປື ຍບ 5670	ອັດຕຣິປື ຍັບ 5670	ອັດຕຣິປື ຍັບ 5670

ຕາຕະລາງທີ 2. ການຝຶກແບບທີ1: train with existing box ໃນTiff/Box 2 ຮູບແບບ

ລຳດັບ [ຮູບປ້າຍລົດ]	Tiff/Box ຮູບແບບທີ1		Tiff/Box ຮູບແບບທີ2	
	Levenshtein Distance	accuracyrate	Levenshtein Distance	accuracyrate
1	6	66.67%	7	61.11%
2	6	66.67%	7	61.11%
3	3	80.00%	6	62.50%
4	6	64.71%	9	52.63%
5	5	72.22%	6	62.50%
6	10	50.00%	13	40.91%
7	9	52.63%	11	42.11%
8	7	61.11%	10	47.37%
9	6	62.50%	8	52.94%
10	8	57.89%	8	52.94%
11	6	66.67%	9	55.00%
12	11	45.00%	13	40.91%
13	8	57.89%	12	42.86%
14	6	64.71%	7	58.82%
15	11	45.00%	12	40.00%
16	5	68.75%	8	52.94%
17	9	43.75%	10	37.50%
18	9	59.09%	12	50.00%
19	9	52.63%	8	52.94%
20	6	66.67%	6	62.50%

ຕາຕະລາງທີ 3. ການຝຶກແບບທີ2: Shape Clustering ໃນTiff/Box 2 ຮູບແບບ

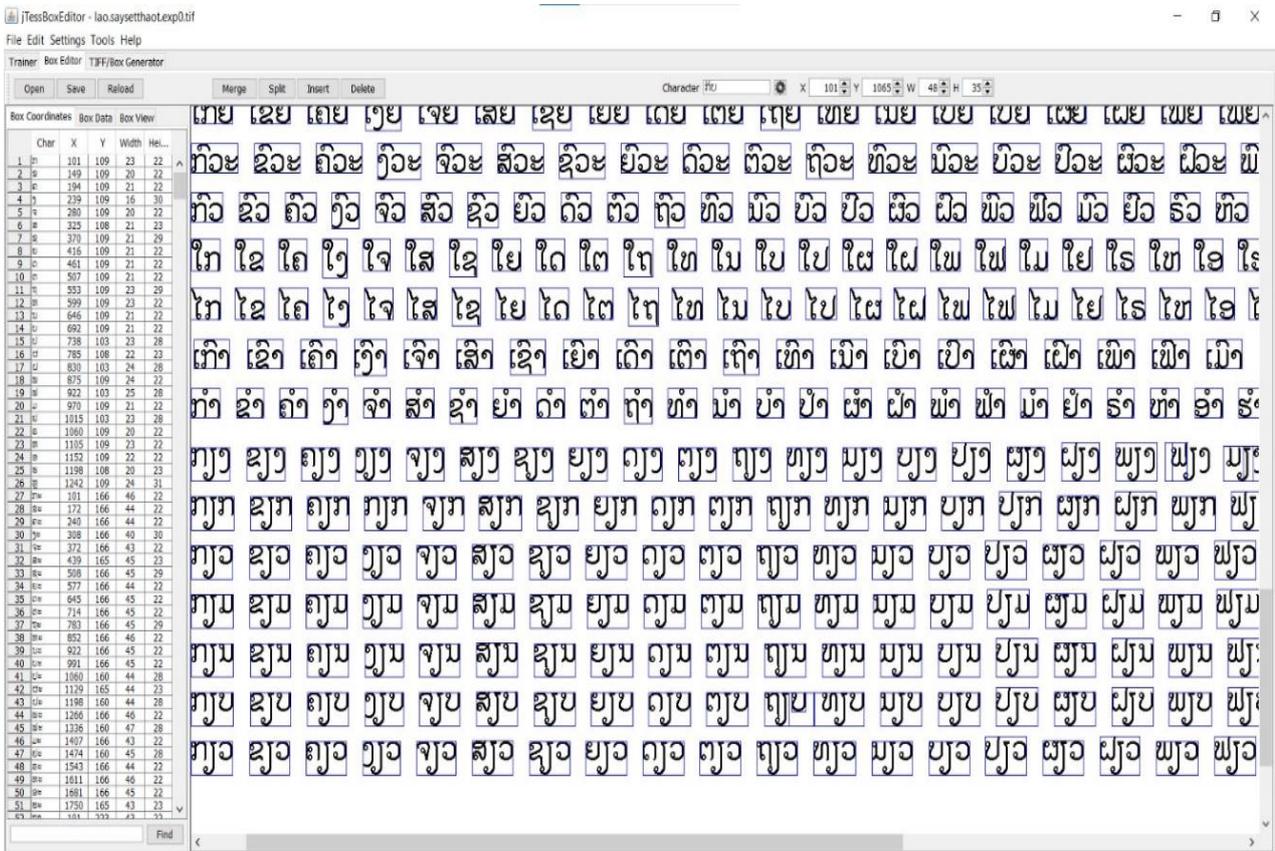
ລຳດັບ [ຮູບປ້າຍ ລິດ]	Tiff/Box ຮູບແບບທີ1		Tiff/Box ຮູບແບບທີ2	
	Levenshtein Distance	accuracyrate	Levenshtein Distance	accuracyrate
1	5	68.75%	5	68.75%
2	5	66.67%	5	66.67%
3	4	73.33%	4	73.33%
4	5	64.29%	5	64.29%
5	5	68.75%	5	68.75%
6	7	58.82%	7	58.82%
7	9	43.75%	8	55.56%
8	7	58.82%	7	58.82%
9	8	55.56%	7	58.82%
10	8	55.56%	7	58.82%
11	6	68.42%	6	68.42%
12	15	28.57%	15	28.57%
13	6	62.50%	6	62.50%
14	5	66.67%	5	66.67%
15	8	52.94%	8	52.94%
16	5	66.67%	5	66.67%
17	9	43.75%	9	43.75%
18	8	63.64%	8	63.64%
19	9	50.00%	9	50.00%
20	6	64.71%	6	64.71%



ຮູບພາບທີ 1: ຕົວຢ່າງ ປ້າຍລິດລາວ

ກຂຄງຈສຊຍດຕຖທນປປຜຜພຟມຢຣຫອຣ  
 ກະ ຂະ ຄະ ງະ ຈະ ສະ ຊະ ຍະ ດະ ຕະ ຖະ ທະ ນະ ປະ ປະ ຜະ ຜະ ພະ ພະ ມະ ຢະ ຣະ ຫະ ອະ ຣະ  
 ກາ ຂາ ຄາ ງາ ຈາ ສາ ຊາ ຍາ ດາ ຕາ ຖາ ທາ ນາ ປາ ຜາ ຜາ ພາ ຟາ ມາ ຢາ ຣາ ຫາ ອາ ຣາ  
 ກິ ຂິ ຄິ ງິ ຈິ ສິ ຊິ ຍິ ດິ ຕິ ຖິ ທິ ນິ ປິ ຜິ ຜິ ພິ ຟິ ມິ ຢິ ຣິ ຫິ ອິ ຣິ  
 ກີ ຂີ ຄີ ງີ ຈີ ສີ ຊີ ຍີ ດີ ຕີ ຖີ ທີ ນີ ປີ ຜີ ຜີ ພີ ຟີ ມີ ຢີ ຣີ ຫີ ອີ ຣີ  
 ກື ຂື ຄື ງື ຈື ສື ຊື ຍື ດື ຕື ຖື ທື ນື ປື ຜື ຜື ພື ຟື ມື ຢື ຣື ຫື ອື ຣື  
 ກຸ ຂຸ ຄຸ ງຸ ຈຸ ສຸ ຊຸ ຍຸ ດຸ ຕຸ ຖຸ ທຸ ນຸ ປຸ ຜຸ ຜຸ ພຸ ຟຸ ມຸ ຢຸ ຣຸ ຫຸ ອຸ ຣຸ  
 ກູ ຂູ ຄູ ງູ ຈູ ສູ ຊູ ຍູ ດູ ຕູ ຖູ ທູ ນູ ປູ ຜູ ຜູ ພູ ຟູ ມູ ຢູ ຣູ ຫູ ອູ ຣູ  
 ດາະ ເຂະ ເຄະ ເງະ ເຈະ ເສະ ເຊະ ເຍະ ເດະ ເຕະ ເຖະ ເທະ ເນະ ເບະ ເວະ ເຜະ ເຜະ ເພະ ເພະ ເມະ ເຢະ ເຣະ ເຫະ ເອະ ເຣະ  
 ດາ ເຂ ເຄ ເງ ເຈ ເສ ເຊ ເຍ ເດ ເຕ ເຖ ເທ ເນ ເບ ເວ ເຜ ເຜ ເພ ເພ ເມ ເຢ ເຣ ເຫ ເອ ເຣ  
 ແກະ ແຂະ ແຄະ ແງະ ແຈະ ແສະ ແຊະ ແຍະ ແດະ ແຕະ ແຖະ ແທະ ແນະ ແບະ ແວະ ແຜະ ແຜະ ແພະ ແພະ ແມະ ແຢະ ແຣະ ແຫະ ແອະ ແຣະ  
 ແກ ແຂ ແຄ ແງ ແຈ ແສ ແຊ ແຍ ແດ ແຕ ແຖ ແທ ແນ ແບ ແວ ແຜ ແຜ ແພ ແພ ແມ ແຢ ແຣ ແຫ ແອ ແຣ  
 ໂກະ ໂຂະ ໂຄະ ໂງະ ໂຈະ ໂສະ ໂຊະ ໂຍະ ໂດະ ໂຕະ ໂຖະ ໂທະ ໂນະ ໂບະ ໂວະ ໂຜະ ໂຜະ ໂພະ ໂພະ ໂມະ ໂຍະ ໂຣະ ໂຫະ ໂອະ ໂຣະ  
 ໂກ ໂຂ ໂຄ ໂງ ໂຈ ໂສ ໂຊ ໂຍ ໂດ ໂຕ ໂຖ ໂທ ໂນ ໂບ ໂວ ໂຜ ໂຜ ໂພ ໂພ ໂມ ໂຍ ໂຣ ໂຫ ໂອ ໂຣ  
 ດາະ ເຂາະ ເຄາະ ເງາະ ເຈາະ ເສາະ ເຊາະ ເຍາະ ເດາະ ເຕາະ ເຖາະ ເທາະ ເນາະ ເບາະ ເວາະ ເຜາະ ເຜາະ ເພາະ ເພາະ ເມາະ ເຢາະ ເຣາະ ເຫາະ ເອາະ ເຣາະ  
 ກໍ ຂໍ ຄໍ ງໍ ຈໍ ສໍ ຊໍ ຍໍ ດໍ ຕໍ ຖໍ ທໍ ນໍ ປໍ ຜໍ ຜໍ ພໍ ຟໍ ມໍ ຢໍ ຣໍ ຫໍ ອໍ ຣໍ  
 ເກີ ເຂີ ເຄີ ເງີ ເຈີ ເສີ ເຊີ ເຍີ ເດີ ເຕີ ເຖີ ເທີ ເນີ ເບີ ເວີ ເຜີ ເຜີ ເພີ ເພີ ເມີ ເຢີ ເຣີ ເຫີ ເອີ ເຣີ  
 ເກີ ເຂີ ເຄີ ເງີ ເຈີ ເສີ ເຊີ ເຍີ ເດີ ເຕີ ເຖີ ເທີ ເນີ ເບີ ເວີ ເຜີ ເຜີ ເພີ ເພີ ເມີ ເຢີ ເຣີ ເຫີ ເອີ ເຣີ  
 ເດັຍ ເຂັຍ ເຄັຍ ເງັຍ ເຈັຍ ເສັຍ ເຊັຍ ເຍັຍ ເດັຍ ເຕັຍ ເຖັຍ ເທັຍ ເນັຍ ເບັຍ ເວັຍ ເຜັຍ ເພັຍ ເພັຍ ເມັຍ ເຢັຍ ເຣັຍ ເຫັຍ ເອັຍ ເຣັຍ  
 ເດຍ ເຂຍ ເຄຍ ເງຍ ເຈຍ ເສຍ ເຊຍ ເຍຍ ເດຍ ເຕຍ ເຖຍ ເທຍ ເນຍ ເບຍ ເວຍ ເຜຍ ເພຍ ເພຍ ເມຍ ເຢຍ ເຣຍ ເຫຍ ເອຍ ເຣຍ  
 ເກືອ ເຂືອ ເຄືອ ເງືອ ເຈືອ ເສືອ ເຊືອ ເຍືອ ເດືອ ເຕືອ ເຖືອ ເທືອ ເນືອ ເບືອ ເວືອ ເຜືອ ເພືອ ເພືອ ເມືອ ເຢືອ ເຣືອ ເຫືອ ເອືອ ເຣືອ  
 ເກື້ຍ ເຂື້ຍ ເຄື້ຍ ເງື້ຍ ເຈື້ຍ ເສື້ຍ ເຊື້ຍ ເຍື້ຍ ເດື້ຍ ເຕື້ຍ ເຖື້ຍ ເທື້ຍ ເນື້ຍ ເບື້ຍ ເວື້ຍ ເຜື້ຍ ເພື້ຍ ເພື້ຍ ເມື້ຍ ເຢື້ຍ ເຣື້ຍ ເຫື້ຍ ເອື້ຍ ເຣື້ຍ  
 ກົວະ ຂົວະ ຄົວະ ງົວະ ຈົວະ ສົວະ ຊົວະ ຍົວະ ດົວະ ຕົວະ ຖົວະ ທົວະ ນົວະ ປົວະ ປົວະ ຜົວະ ຜົວະ ພົວະ ພົວະ ມົວະ ມົວະ ອົວະ ອົວະ  
 ກົວ ຂົວ ຄົວ ງົວ ຈົວ ສົວ ຊົວ ຍົວ ດົວ ຕົວ ຖົວ ທົວ ນົວ ປົວ ປົວ ຜົວ ຜົວ ພົວ ພົວ ມົວ ຍົວ ອົວ ອົວ ອົວ  
 ໂກ ໂຂ ໂຄ ໂງ ໂຈ ໂສ ໂຊ ໂຍ ໂດ ໂຕ ໂຖ ໂທ ໂນ ໂບ ໂວ ໂຜ ໂຜ ໂພ ໂພ ໂມ ໂຍ ໂຣ ໂຫ ໂອ ໂຣ  
 ໂກ ໂຂ ໂຄ ໂງ ໂຈ ໂສ ໂຊ ໂຍ ໂດ ໂຕ ໂຖ ໂທ ໂນ ໂບ ໂວ ໂຜ ໂຜ ໂພ ໂພ ໂມ ໂຍ ໂຣ ໂຫ ໂອ ໂຣ  
 ເກົາ ເຂົາ ເຄົາ ເງົາ ເຈົາ ເສົາ ເຊົາ ເຍົາ ເດົາ ເຕົາ ເຖົາ ເທົາ ເນົາ ເບົາ ເວົາ ເຜົາ ເພົາ ເພົາ ເມົາ ເຢົາ ເຣົາ ເຫົາ ເອົາ ເຣົາ  
 ກໍາ ຂໍາ ຄໍາ ງໍາ ຈໍາ ສໍາ ຊໍາ ຍໍາ ດໍາ ຕໍາ ຖໍາ ທໍາ ນໍາ ປໍາ ຜໍາ ຜໍາ ພໍາ ຟໍາ ມໍາ ຢໍາ ຣໍາ ຫໍາ ອໍາ ຣໍາ  
 0 1 2 3 4 5 6 7 8 9

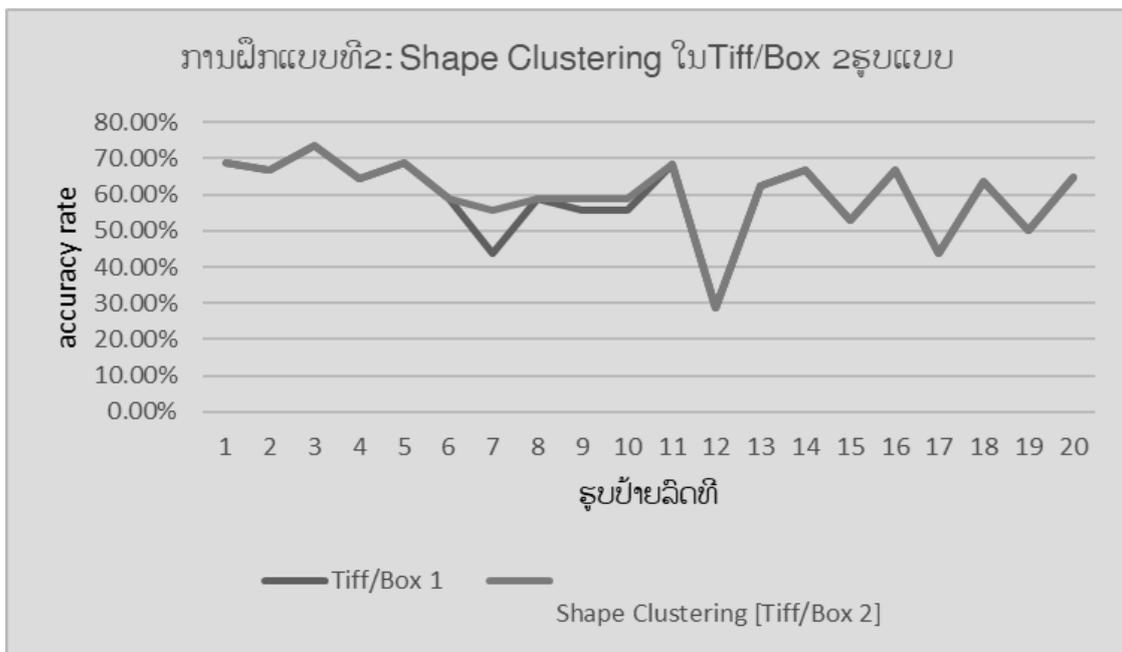
ຮູບພາບທີ 2: ຂັ້ນຕອນການຈັດຈຳຕົວເລກ ແລະ ຕົວອັກສອນ ບ້າຍລົດລາວ



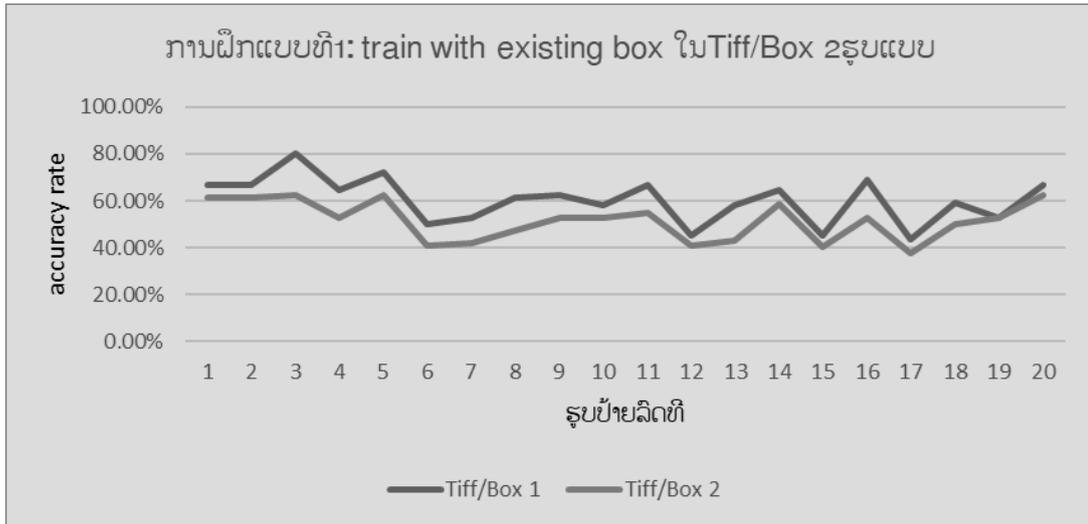
ຮູບພາບທີ 3: ຊຸດຂໍ້ມູນ

1 ເຊກອງ ຍປ 6155	2 ເຊກອງ ຂຜ 5514	3 ຫົວໝັນ ຖຜ 8732	4 ເຊກອງ ພຂ 0072
5 ອັດຕະປື ປຢ 2737	6 ຫຼວງນໍ້າທາ ຕຈ 8099	7 ສາລະວັນ ຢດ 8417	8 ສາລະວັນ ຄພ 2054
9 ບໍ່ແກ້ວ ບຍ 1446	10 ສາລະວັນ ຂຮ 0042	11 ໄຊສົມບູນ ປຄ 0047	12 ຜົ້ງສາລີ ຜຜ 0044
13 ໄຊຍະບູລີ ຜປ 7168	14 ຫົວໝັນ ຂທ 6830	15 ຫຼວງນໍ້າທາ ຈກ 6212	16 ຫົວໝັນ ສຈ 0442
17 ຜົ້ງສາລີ ນຜ 0864	18 ໄຊສົມບູນ ງຫ 1344	19 ສາລະວັນ ຂອ 9345	20 ອັດຕະປື ຍບ 5670

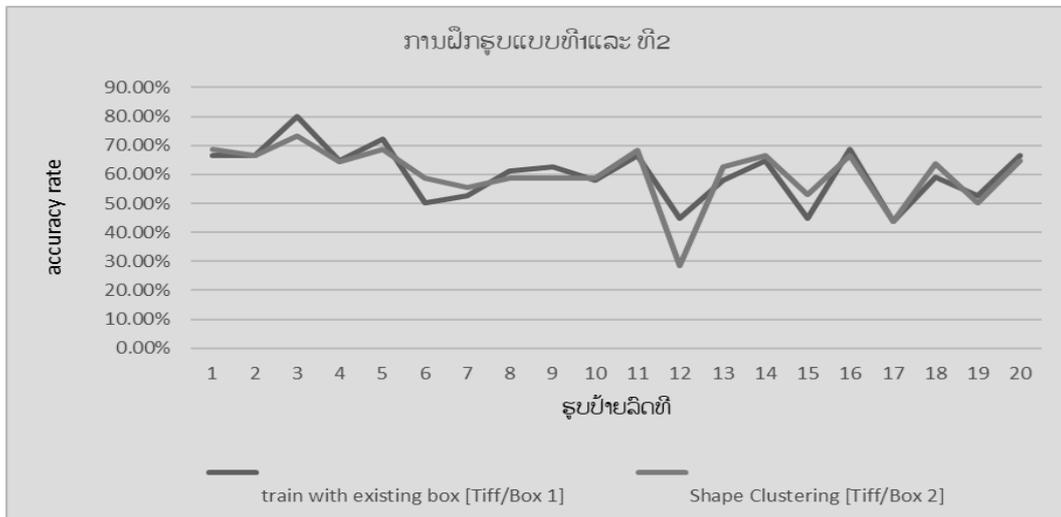
ຮູບພາບທີ 4: ຕົວຢ່າງການສ້າງ Tiff/Box ຕາມຊຸດຂໍ້ມູນ



ຮູບພາບທີ 5: ສະແດງເປີເຊັນຄວາມແມ້ນຍໍາການຝຶກແບບທີ1 train with existing box ໃນTiff/Box 2 ຮູບແບບ



ຮູບພາບທີ 6: ສະແດງເປີເຊັນຄວາມແມ້ນຍຳການຝຶກແບບທີ2 Shape Clustering ໃນTiff/Box 2 ຮູບແບບ



ຮູບພາບທີ 7: ສະແດງເປີເຊັນການປຽບທຽມຄວາມແມ້ນຍຳຂອງທັງສອງຮູບແບບ