



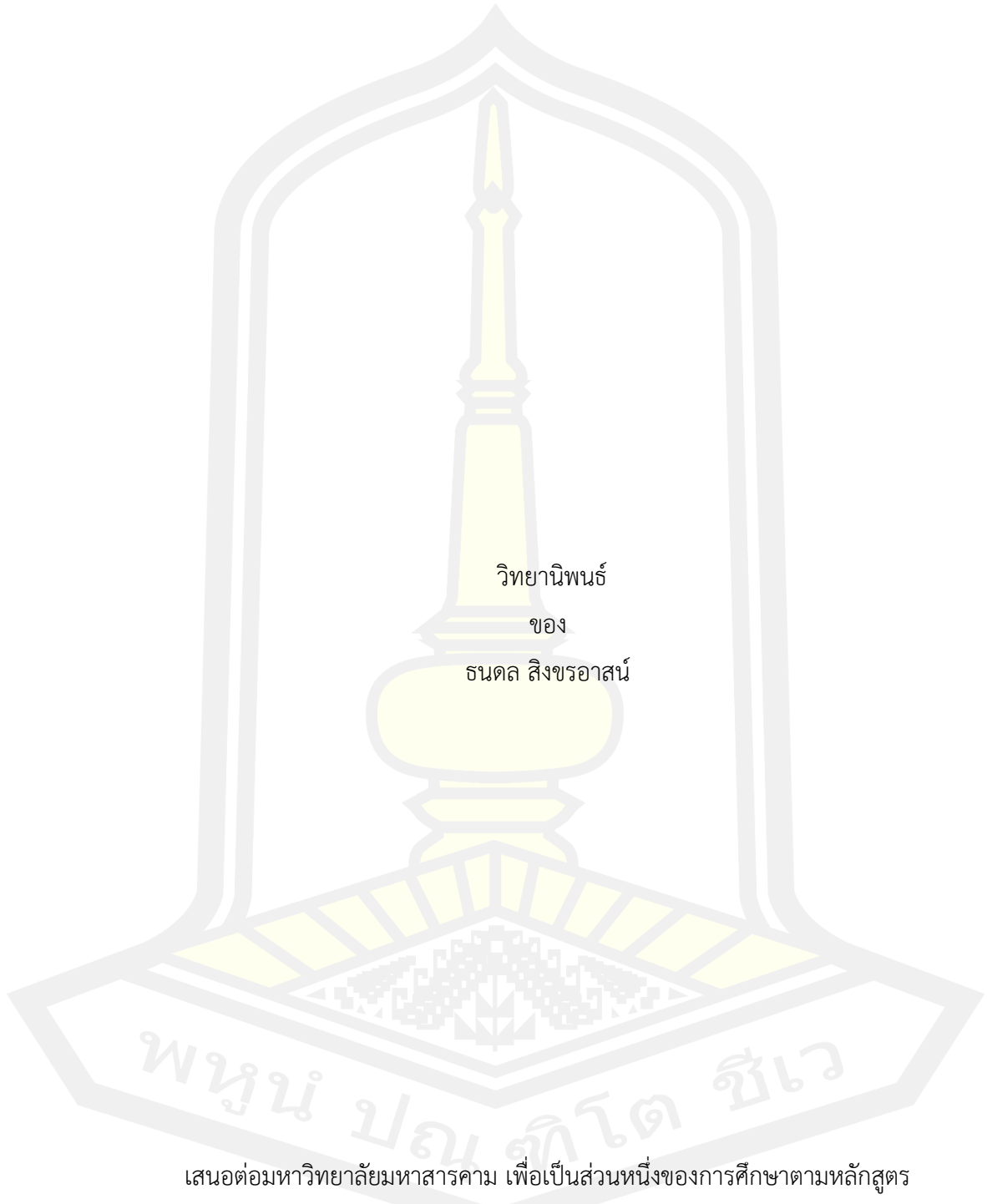
การเรียนรู้เชิงลึกสำหรับการตรวจจับและรู้จำคำบรรยายในวิดีโอ

วิทยานิพนธ์
ของ
ธนดล สิงขรอาสน์

เสนอต่อมหาวิทยาลัยมหาสารคาม เพื่อเป็นส่วนหนึ่งของการศึกษาตามหลักสูตร
ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ
ธันวาคม 2564

ลิขสิทธิ์เป็นของมหาวิทยาลัยมหาสารคาม

การเรียนรู้เชิงลึกสำหรับการตรวจจับและรู้จำคำบรรยายในวิดีโอ



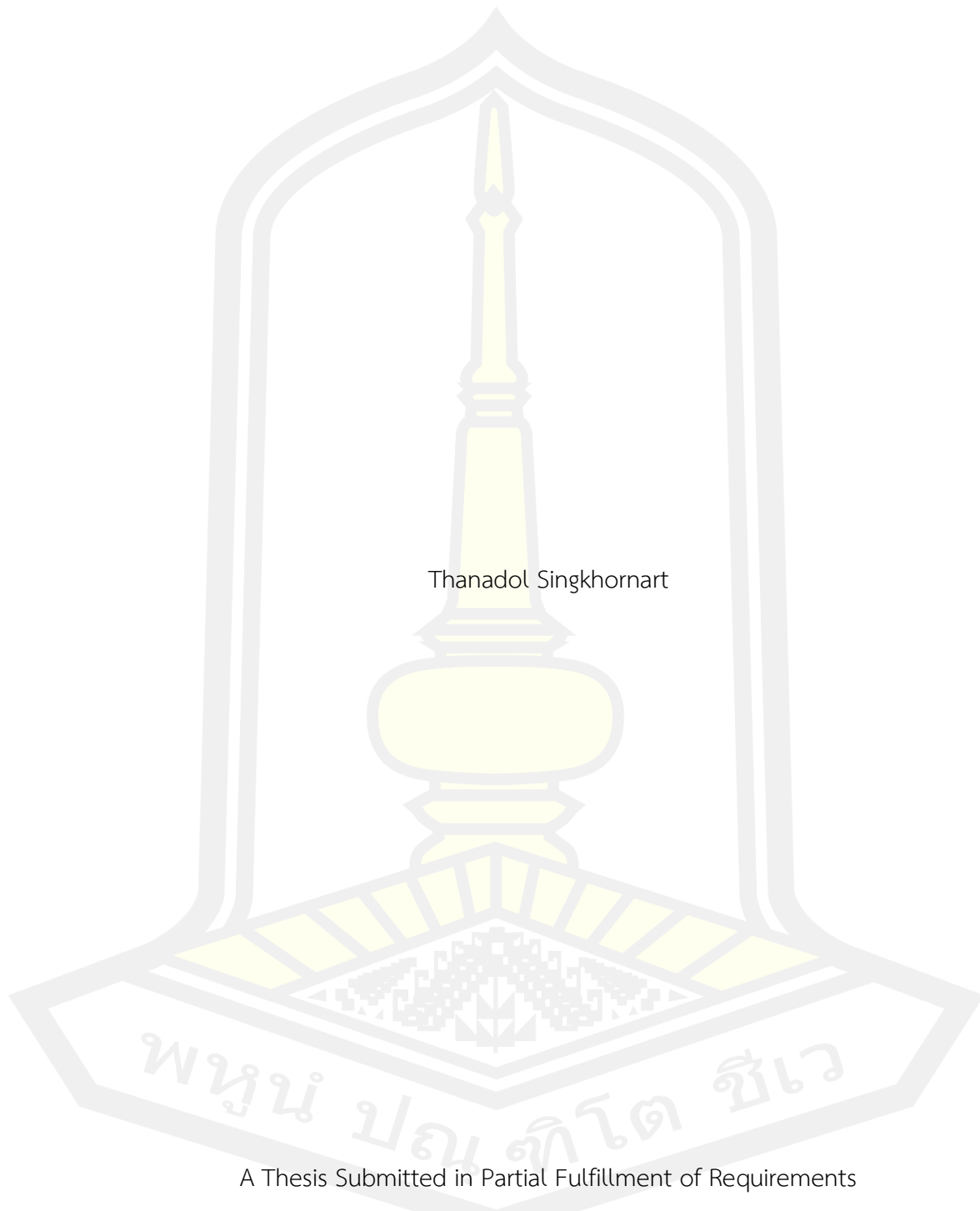
เสนอต่อมหาวิทยาลัยมหาสารคาม เพื่อเป็นส่วนหนึ่งของการศึกษาตามหลักสูตร

ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ

ธันวาคม 2564

ลิขสิทธิ์เป็นของมหาวิทยาลัยมหาสารคาม

Deep Learning for Video Subtitle Detection and Recognition



Thanadol Singhornart

A Thesis Submitted in Partial Fulfillment of Requirements
for Master of Science (Information Technology)

December 2021

Copyright of Mahasarakham University



คณะกรรมการสอบวิทยานิพนธ์ ได้พิจารณาวิทยานิพนธ์ของนายธนดล สิงขรอาสน์
แล้วเห็นสมควรรับเป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต
สาขาวิชาเทคโนโลยีสารสนเทศ ของมหาวิทยาลัยมหาสารคาม

คณะกรรมการสอบวิทยานิพนธ์

.....ประธานกรรมการ

(รศ. ดร. ธัชพงศ์ กัตถัญญกุล)

.....อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก

(ผศ. ดร. โอฬาริก สุรินดียะ)

.....กรรมการ

(ผศ. ดร. สาธิต แสงประดิษฐ์)

.....กรรมการ

(ผศ. ดร. แกมกาญจน์ สมประเสริฐศรี)

มหาวิทยาลัยอนุมัติให้รับวิทยานิพนธ์ฉบับนี้ เป็นส่วนหนึ่งของการศึกษาตามหลักสูตร
ปริญญา วิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ ของมหาวิทยาลัยมหาสารคาม

.....
(ผศ. ศศิธร แก้วมัน)

คณบดีคณะวิทยาการสารสนเทศ

.....
(รศ. ดร. กริสน์ ชัยมูล)

คณบดีบัณฑิตวิทยาลัย

ชื่อเรื่อง	การเรียนรู้เชิงลึกสำหรับการตรวจจับและรู้จำคำบรรยายในวิดีโอ		
ผู้วิจัย	ชนดล สิงขรอาสน์		
อาจารย์ที่ปรึกษา	ผู้ช่วยศาสตราจารย์ ดร. โอฬาริก สุรินตะ		
ปริญญา	วิทยาศาสตรมหาบัณฑิต	สาขาวิชา	เทคโนโลยีสารสนเทศ
มหาวิทยาลัย	มหาวิทยาลัยมหาสารคาม	ปีที่พิมพ์	2564

บทคัดย่อ

ในปัจจุบันมีวิดีโอจำนวนมากที่ถูกเผยแพร่ผ่านอินเทอร์เน็ตในช่องทางต่าง ๆ เช่น Youtube และ Facebook มีผู้ชมบางส่วนที่มีปัญหาในการรับรู้ข้อมูลจากวิดีโอเนื่องจากปัญหาทางด้านภาษาหรือมีปัญหาด้านการฟัง ดังนั้นคำบรรยายจึงถูกเพิ่มเข้ามาในวิดีโอ ในวิทยานิพนธ์นี้ได้นำเสนอถึงการนำวิธีการเรียนรู้เชิงลึกมาใช้โดยใช้วิธีโครงข่ายประสาทแบบคอนโวลูชัน (CNN) ร่วมกับ วิธีหน่วยความจำระยะสั้นระยะยาว (LSTM) ซึ่งเรียกว่า CNN-LSTM เพื่อที่จะนำมารู้จำคำบรรยายจากวิดีโอ เราได้สร้างตัวอย่างต้นแบบ CNN ที่มีจำนวน 16 ชั้น โดยชั้นสุดท้ายเป็น การย่อขนาดโดยใช้ค่าสูงสุด (Max-pooling) ก่อนที่จะส่งเข้า LSTM โดยในการเรียนรู้เราได้ใช้รูปภาพคำบรรยายที่มีรูปแบบ ขนาด และพื้นหลังที่หลากหลาย แล้วใช้ การจำแนกการเชื่อมต่อชั่วคราว (CTC loss) ในการคำนวณค่า loss และถอดรหัสเป็นผลลัพธ์ สำหรับข้อมูลที่เราใช้ในการเรียนรู้ได้มาจากการรวบรวม 24 วิดีโอจาก Youtube และ Facebook ที่มีคำบรรยายภาษาไทย อังกฤษ ตัวเลขไทยและตัวเลขอารบิก ซึ่งมีทั้งหมด 157 ตัวเพื่อนำมาถอดรหัสข้อมูลในชุดรูปภาพนั้นมีทั้งหมด 4,224 รูป ซึ่งได้ค่าเฉลี่ยความผิดพลาดที่น้อยที่สุดคือ 11.06%

คำสำคัญ : การรู้จำคำบรรยายวิดีโอ, โครงข่ายประสาทคอนโวลูชัน, หน่วยความจำระยะสั้นระยะยาว, การจำแนกการเชื่อมต่อชั่วคราว

พนุณ ปณฺ ทิโต ชิว

TITLE	Deep Learning for Video Subtitle Detection and Recognition		
AUTHOR	Thanadol Singkhornart		
ADVISORS	Assistant Professor Olarik Surinta , Ph.D.		
DEGREE	Master of Science	MAJOR	Information Technology
UNIVERSITY	Maharakham University	YEAR	2021

ABSTRACT

Nowadays, many videos have been published on Internet channels such as Youtube and Facebook. Many audiences, however, cannot understand the contents of the video, maybe due to the different languages and even hearing impairment. As a result, subtitles have been added to videos. In this paper, we proposed deep learning techniques, which were the combination between convolutional neural networks (CNN) and long short-term memory (LSTM) networks, called CNN-LSTM, to recognize video subtitles. We created the simplified CNN architecture with 16 weight layers. The last layer of the CNN was downsampling using max-pooling before sending it to the LSTM network. We first trained our CNN-LSTM architecture on printed text data which contained various font styles, diverse font sizes, and complicated backgrounds. The connectionist temporal classification was then used as a loss function to calculate the loss value and decode the output of the network. For the video subtitle dataset, we collected 24 videos from Youtube and Facebook, consisting of Thai, English, Arabic, and Thai numbers. The dataset also contained 157 characters. In this dataset, we extracted 4,224 subtitle images from the videos. The proposed CNN-LSTM architecture achieved an average character error rate of 11.06%.

Keyword : Video subtitle text recognition, Convolutional neural networks, long short-term memory network, Connectionist temporal classification

กิตติกรรมประกาศ

วิทยานิพนธ์นี้สำเร็จได้ด้วยความกรุณาของอาจารย์โอฬาริก สุรินตะ อาจารย์ที่ปรึกษา วิทยานิพนธ์ที่ได้ให้คำแนะนำ แนวคิดตลอดจนแก้ไขข้อบกพร่องต่าง ๆ มาโดยตลอด จนวิทยานิพนธ์เล่มนี้เสร็จสมบูรณ์ ขอกราบขอบพระคุณ ณ ที่นี้

ธนดล สิงขรอาสน์



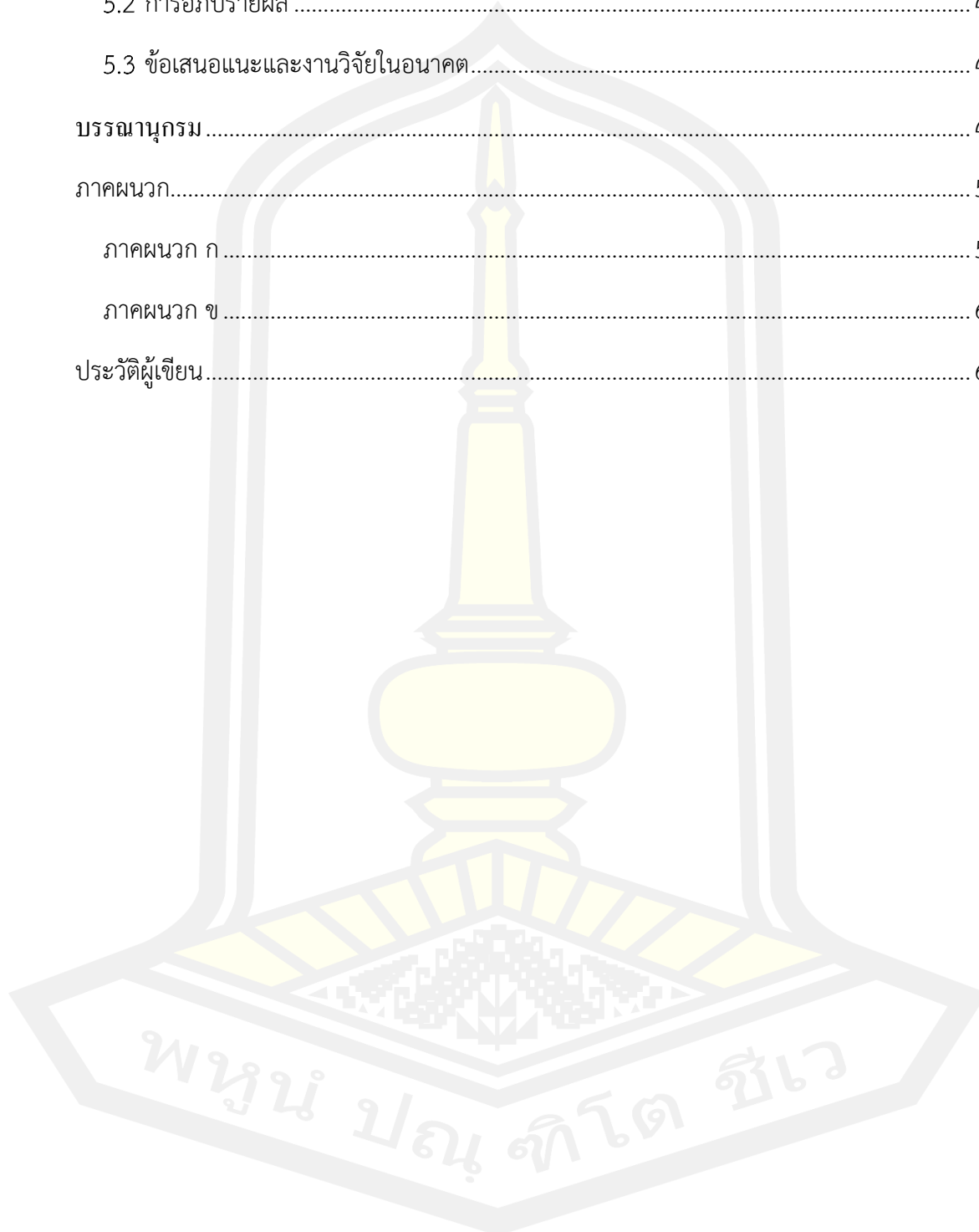
สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	ง
บทคัดย่อภาษาอังกฤษ.....	จ
กิตติกรรมประกาศ.....	ฉ
สารบัญ.....	ช
สารบัญภาพ	ฉ
สารบัญตาราง.....	ท
บทที่ 1 บทนำ	1
1.1 ความเป็นมาของการวิจัย	1
1.2 ความมุ่งหมายของการวิจัย	3
1.3 ขอบเขตของการวิจัย.....	3
1.3.1 ข้อมูลที่ใช้ในการวิจัย.....	3
1.3.1.1 รวบรวมวิดิทัศน์จำนวน 24 วิดิทัศน์	3
1.3.1.2 เตรียมข้อมูลเพื่อนำไปใช้ในกระบวนการตรวจจับคำบรรยายวิดิทัศน์.....	4
1.3.1.3 เตรียมข้อมูลเพื่อนำไปใช้ในกระบวนการรู้จำคำบรรยายวิดิทัศน์	4
1.3.2 ตัวอักษรที่ใช้ในการรู้จำคำบรรยายวิดิทัศน์.....	5
1.3.3 แบ่งข้อมูลที่ใช้ในการทดลองและวิธีการวัดประสิทธิภาพ	5
1.3.3.1 ข้อมูลที่ใช้ในการตรวจจับคำบรรยายวิดิทัศน์.....	5
1.3.3.2 ข้อมูลที่ใช้ในการรู้จำคำบรรยายวิดิทัศน์	6
1.3.4 กระบวนการเรียนรู้เชิงลึกสำหรับการตรวจจับและรู้จำคำบรรยายวิดิทัศน์.....	6
1.3.4.1 ตรวจจับคำบรรยายวิดิทัศน์ ด้วยวิธีการเรียนรู้เชิงลึก.....	6
1.3.4.2 รู้จำคำบรรยายวิดิทัศน์ ด้วยวิธีการเรียนรู้เชิงลึก	6

1.4 ผลที่คาดว่าจะได้รับจากงานวิจัยครั้งนี้	6
1.5 นิยามศัพท์เฉพาะ.....	6
บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง	7
2.1 การเรียนรู้เชิงลึก (Deep Learning).....	7
2.2 โครงข่ายประสาทคอนโวลูชัน (Convolution Neural Network: CNN).....	7
2.2.1 ชั้นรับข้อมูล (Input Layer)	8
2.2.2 ชั้นคอนโวลูชัน (Convolution Layer).....	9
2.2.3 ชั้นพูลลิ่ง (Pooling)	10
2.2.4 ชั้นเชื่อมโยงสมบูรณ์ (Fully Connected Layer).....	11
2.2.5 ชั้นผลลัพธ์ (output).....	11
2.3 การตรวจจับคำบรรยาย (Subtitle Detection)	12
2.3.1 YOLOv3.....	12
2.3.2 Tiny-YOLOv3.....	14
2.3.3 RetinaNet.....	15
2.4 การรู้จำคำบรรยาย (Subtitle Recognition).....	15
2.4.1 Visual Geometry Group (VGG).....	16
2.4.2 Long Short-Term Memory (LSTM).....	17
2.4.3 Gated Recurrent Unit (GRU).....	18
2.4.4 Connectionist Temporal Classification Loss (CTC Loss)	18
2.5 การตรวจสอบประสิทธิภาพ.....	19
2.5.1 Intersection Over Union (IoU).....	19
2.5.2 อัตราความแม่นยำเฉลี่ย (mean Average Precision: mAP)	20
2.5.2.1 Precision	20
2.5.2.2 Recall	20

2.5.3 ความผิดพลาดระดับตัวอักษร (Character Error Rate: CER).....	20
2.6 งานวิจัยที่เกี่ยวข้อง	20
2.6.1 งานวิจัยที่เกี่ยวข้องกับการตรวจจับคำบรรยาย (Text Detection)	20
2.6.2 งานวิจัยที่เกี่ยวข้องกับการรู้จำคำบรรยาย (Text Recognition)	22
บทที่ 3 วิธีดำเนินการวิจัย	24
3.1 ชุดข้อมูลที่ใช้ในงานวิจัย.....	24
3.1.1 รูปภาพที่ใช้สำหรับตรวจจับคำบรรยายวีดิทัศน์	25
3.1.2 รูปภาพคำบรรยายวีดิทัศน์ที่ใช้สำหรับการรู้จำ	27
3.2 การเตรียมข้อมูล.....	28
3.2.1 หาตัวอย่างที่มีการฝังคำบรรยายภาษาไทยและอังกฤษ	28
3.2.2 นำคลิปวีดิทัศน์ไปแปลงเป็นรูปภาพ.....	28
3.2.3 นำรูปภาพไปสร้าง ground truth.....	29
3.2.4 นำรูปภาพคำบรรยาย (Subtitle image) มาสร้างไฟล์ข้อความตัวอักษรประกอบ (Text transcription).....	30
3.3 การตรวจจับคำบรรยายวีดิทัศน์.....	30
3.3.1 ข้อมูลที่ใช้ในการทดสอบกระบวนการตรวจจับคำบรรยายวีดิทัศน์.....	30
3.3.2 ทำการสร้างโมเดลในการตรวจจับคำบรรยายวีดิทัศน์	30
3.3.3 ทดสอบประสิทธิภาพของโมเดลในการตรวจจับคำบรรยายวีดิทัศน์	31
3.4 การรู้จำคำบรรยายวีดิทัศน์	31
3.4.1 ข้อมูลที่ใช้ในการทดสอบในการรู้จำคำบรรยายวีดิทัศน์.....	31
3.4.2 ทำการสร้างโมเดลในการรู้จำคำบรรยายวีดิทัศน์	31
3.4.4 ทดสอบประสิทธิภาพของโมเดลในการรู้จำคำบรรยายวีดิทัศน์.....	34
บทที่ 4 ผลลัพธ์การทดลอง.....	36
บทที่ 5 สรุปและข้อเสนอแนะ	43

5.1 สรุปผลการวิจัย.....	43
5.2 การอภิปรายผล	43
5.3 ข้อเสนอแนะและงานวิจัยในอนาคต.....	44
บรรณานุกรม	46
ภาคผนวก.....	51
ภาคผนวก ก	52
ภาคผนวก ข	61
ประวัติผู้เขียน	68



สารบัญภาพ

ภาพประกอบที่ 1.1 ตัวอย่างของคำบรรยายที่ปรากฏในวิดีโอที่คน ก) คำบรรยายที่ปรากฏบริเวณ ด้านล่างของวิดีโอ และ ข) คำบรรยายที่ปรากฏขึ้นในบริเวณที่น่าสนใจ	1
ภาพประกอบที่ 1.2 ตัวอย่างของการรู้จำคำบรรยาย ก) ภาพที่ได้จากขั้นตอนการตรวจจับคำบรรยาย และ ข) ผลลัพธ์ที่ได้จากการรู้จำคำบรรยาย โดยผลลัพธ์จะอยู่ในรูปแบบของ ข้อความ (Text)	2
ภาพประกอบที่ 1.3 ตัวอย่างของคำบรรยายวิดีโอที่ใช้ในวิทยานิพนธ์	3
ภาพประกอบที่ 1.4 ตัวอย่างการสร้าง Ground Truth ในบริเวณที่เป็นคำบรรยายวิดีโอ.....	4
ภาพประกอบที่ 1.5 ตัวอย่างการเก็บข้อมูลคำบรรยายวิดีโอ 964_ตัวระบุเปิดเคลย์มอร์ หรือ รุ่น ระบุเปิดสังหารบุคคล [964 คือหมายเลขลำดับรูปภาพ และ “ตัวระบุเปิดเคลย์มอร์ หรือ รุ่นระบุเปิดสังหารบุคคล” คือคำตอบ (Label)].....	4
ภาพประกอบที่ 2.1 การเรียนรู้เชิงลึก (Deep Learning) [10]	7
ภาพประกอบที่ 2.2 สถาปัตยกรรมโครงข่ายประสาทคอนโวลูชัน (Convolution Neural Network: CNN) [11]	8
ภาพประกอบที่ 2.3 ตัวอย่างการทำงานของชั้นคอนโวลูชัน ขนาด 3x3 stride=1 [11].....	9
ภาพประกอบที่ 2.4 ตัวอย่างการทำงานของชั้น Max Pooling [11]	10
ภาพประกอบที่ 2.5 ตัวอย่างการเชื่อมต่อระหว่างชั้นรับข้อมูล ชั้นเชื่อมโยงสมบูรณ์	11
ภาพประกอบที่ 2.6 สถาปัตยกรรม Darknet-53 [3].....	13
ภาพประกอบที่ 2.7 สถาปัตยกรรม Tiny-YOLOv3 [14].....	14
ภาพประกอบที่ 2.8 สถาปัตยกรรม RetinaNet [5].....	15
ภาพประกอบที่ 2.9 สถาปัตยกรรม VGG16 [18]	16
ภาพประกอบที่ 2.10 การทำงานของ LSTM.....	17
ภาพประกอบที่ 2.11 การทำงานของ GRU	18

ภาพประกอบที่ 2.12 Area of Overlap โดยสีฟ้าคือพื้นที่ Ground Truth สีแดงคือผลลัพธ์ที่ได้จากการตรวจจับตำแหน่งคำบรรยาย สีเหลืองคือพื้นที่ซ้อนทับ (Area of Overlap)	19
ภาพประกอบที่ 2.13 Area of Union ซึ่งเป็นการรวมของพื้นที่ Ground Truth และผลลัพธ์การตรวจจับตำแหน่งคำบรรยาย	19
ภาพประกอบที่ 2.14 ค่า mAP 62.912 ซึ่งมากกว่ามีค่า IoU มากกว่าร้อยละ 50	19
ภาพประกอบที่ 2.1 การเรียนรู้เชิงลึก (Deep Learning) [10]	7
ภาพประกอบที่ 2.2 สถาปัตยกรรมโครงข่ายประสาทคอนโวลูชัน (Convolution Neural Network: CNN) [11]	8
ภาพประกอบที่ 2.3 ตัวอย่างการทำงานของชั้นคอนโวลูชัน ขนาด 3x3 stride=1 [11]	9
ภาพประกอบที่ 2.4 ตัวอย่างการทำงานของชั้น Max Pooling [11]	10
ภาพประกอบที่ 2.5 ตัวอย่างการเชื่อมต่อระหว่างชั้นรับข้อมูล ชั้นเชื่อมโยงสมบูรณ์	11
ภาพประกอบที่ 2.6 สถาปัตยกรรม Darknet-53 [3]	13
ภาพประกอบที่ 2.7 สถาปัตยกรรม Tiny-YOLOv3 [14]	14
ภาพประกอบที่ 2.8 สถาปัตยกรรม RetinaNet [5]	15
ภาพประกอบที่ 2.9 สถาปัตยกรรม VGG16 [18]	16
ภาพประกอบที่ 2.10 การทำงานของ LSTM	17
ภาพประกอบที่ 2.11 การทำงานของ GRU	18
ภาพประกอบที่ 2.12 Area of Overlap โดยสีฟ้าคือพื้นที่ Ground Truth สีแดงคือผลลัพธ์ที่ได้จากการตรวจจับตำแหน่งคำบรรยาย สีเหลืองคือพื้นที่ซ้อนทับ (Area of Overlap)	19
ภาพประกอบที่ 2.13 Area of Union ซึ่งเป็นการรวมของพื้นที่ Ground Truth และผลลัพธ์การตรวจจับตำแหน่งคำบรรยาย	19
ภาพประกอบที่ 2.14 ค่า mAP 62.912 ซึ่งมากกว่ามีค่า IoU มากกว่าร้อยละ 50	19
ภาพประกอบที่ 4.1 ภาพผลลัพธ์การตรวจจับตำแหน่งคำบรรยาย รวมไปถึงร้อยละที่ถูกต้องและตำแหน่งของกรอบที่ถูกต้องตามภาพ ก) ภาพผลลัพธ์ที่มีกรอบมากกว่า ground	

truth, ข) ภาพผลลัพธ์ที่มีกรอบเท่ากับ ground truth และ ค) ภาพผลลัพธ์ที่มี
 กรอบน้อยกว่า ground truth 37

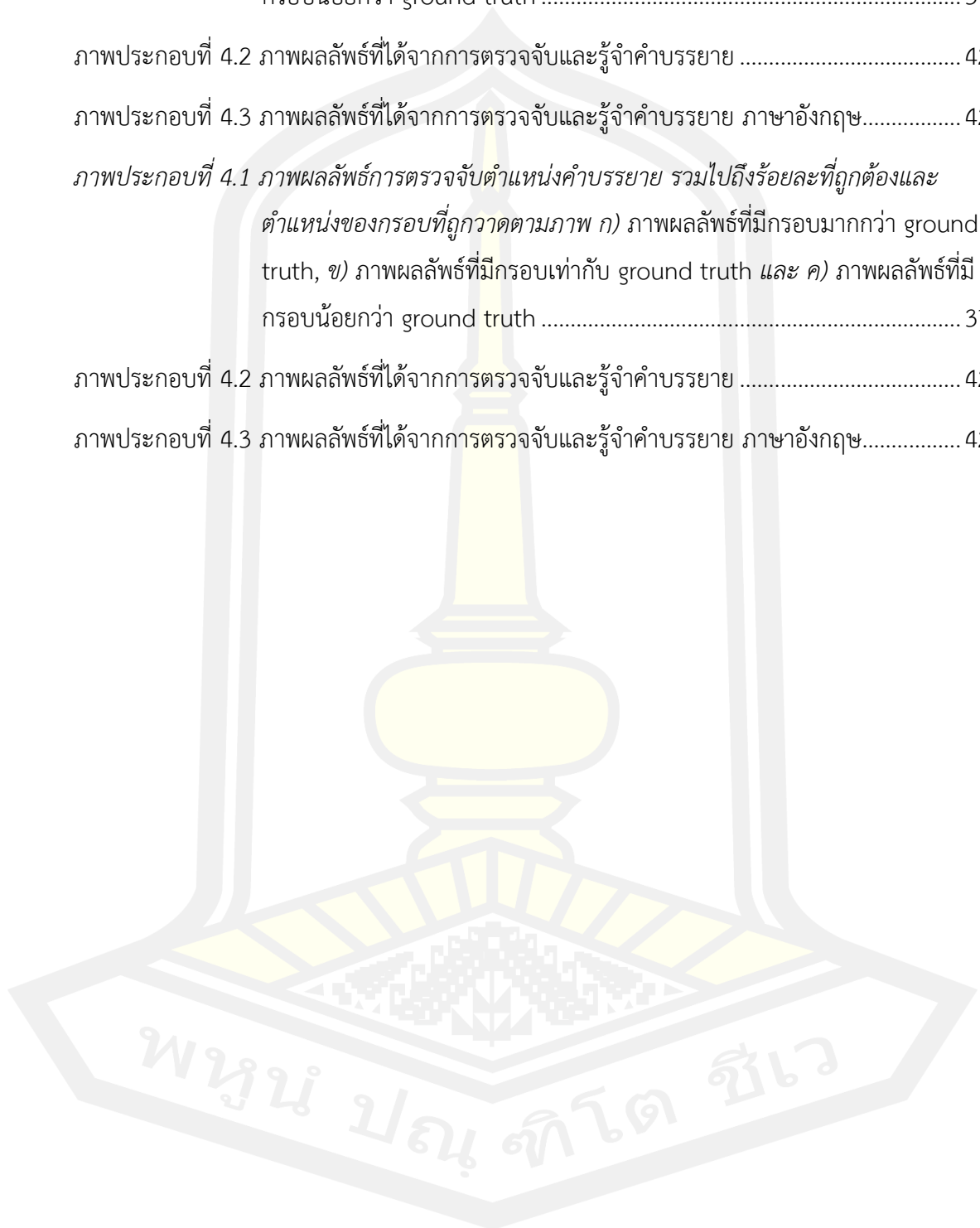
ภาพประกอบที่ 4.2 ภาพผลลัพธ์ที่ได้จากการตรวจจับและรู้จำคำบรรยาย 42

ภาพประกอบที่ 4.3 ภาพผลลัพธ์ที่ได้จากการตรวจจับและรู้จำคำบรรยาย ภาษาอังกฤษ..... 42

ภาพประกอบที่ 4.1 ภาพผลลัพธ์การตรวจจับตำแหน่งคำบรรยาย รวมไปถึงร้อยละที่ถูกต้องและ
 ตำแหน่งของกรอบที่ถูกต้องตามภาพ ก) ภาพผลลัพธ์ที่มีกรอบมากกว่า ground
 truth, ข) ภาพผลลัพธ์ที่มีกรอบเท่ากับ ground truth และ ค) ภาพผลลัพธ์ที่มี
 กรอบน้อยกว่า ground truth 37

ภาพประกอบที่ 4.2 ภาพผลลัพธ์ที่ได้จากการตรวจจับและรู้จำคำบรรยาย 42

ภาพประกอบที่ 4.3 ภาพผลลัพธ์ที่ได้จากการตรวจจับและรู้จำคำบรรยาย ภาษาอังกฤษ..... 42



สารบัญตาราง

ตารางที่ 1.1 ตัวอักษรที่ใช้ในการรู้จำคำบรรยายวิดีโอ.....	5
ตารางที่ 3.1 ตัวอย่างสถาปัตยกรรม CNN-LSTM	32
ตารางที่ 3.2 สถาปัตยกรรมที่ใช้ในการสร้างแบบจำลองการรู้จำคำบรรยายวิดีโอ.....	33
ตารางที่ 3.3 สถาปัตยกรรมตามบทความของ Chamchong et al.	34
ตารางที่ 3.1 ตัวอย่างสถาปัตยกรรม CNN-LSTM.....	32
ตารางที่ 3.2 สถาปัตยกรรมที่ใช้ในการสร้างแบบจำลองการรู้จำคำบรรยายวิดีโอ.....	33
ตารางที่ 3.3 สถาปัตยกรรมตามบทความของ Chamchong et al.	34
ตารางที่ 4.1 ผลลัพธ์ค่า mAP โดยกำหนดให้ค่า IoU = 0.5.....	36
ตารางที่ 4.2 ค่า CER ที่ได้จากสถาปัตยกรรมที่ใช้วิธี CNN และ LSTM	38
ตารางที่ 4.3 ค่า CER ที่ได้จากการใช้สถาปัตยกรรมตามบทความของ Chamchong et al.	38
ตารางที่ 4.4 ตารางเปรียบเทียบค่า CER และ เวลาในการเรียนรู้ (Train).....	39
ตารางที่ 4.5 ผลลัพธ์การทดลองการรู้จำคำบรรยาย.....	39

บทที่ 1

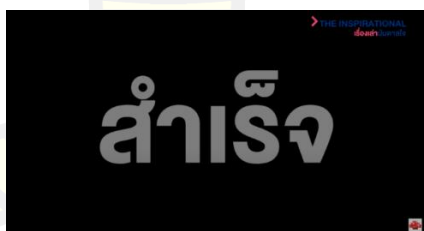
บทนำ

1.1 ความเป็นมาของการวิจัย

สื่อประเภทวีดิทัศน์ได้ถูกเผยแพร่ในหลายช่องทางเช่น ยูทูบ (Youtube) เฟซบุ๊ก (Facebook) และอินสตาแกรม (Instagram) เป็นต้น ทำให้ผู้ชมสามารถเลือกชมได้อย่างอิสระ ทั้งนี้เพื่อให้เข้าถึงผู้ชมได้หลากหลาย เช่น ชาวต่างชาติ และผู้พิการทางการได้ยิน ทำให้ผู้ผลิตวีดิทัศน์เพิ่มคำบรรยาย (Subtitle) ลงไปในวีดิทัศน์ ทำให้เกิดประโยชน์ต่อผู้ชมที่จะได้รับรู้เนื้อหาที่ถูกต้อง และยังเป็นการเพิ่มจำนวนคนเข้าชมวีดิทัศน์ ทั้งนี้ยังส่งผลให้ผู้ผลิตได้รับรายได้จากการเข้าชมวีดิทัศน์ ดังนั้นการเพิ่มคำบรรยายที่อยู่ในวีดิทัศน์สามารถเป็นได้ทั้งคำบรรยายที่ที่แสดงออกมาตามคำพูดของผู้พูดในวีดิทัศน์และ เป็นคำบรรยายที่ปรากฏขึ้นในจุดที่น่าสนใจ แสดงดังภาพประกอบที่ 1.1 ทั้งนี้ ขนาดของตัวอักษรอาจมีความแตกต่างกันออกไป เพื่อให้ผู้ชมเกิดความสนใจและเข้าชมวีดิทัศน์อื่น ๆ ต่อไป



ก)



ข)

ภาพประกอบที่ 1.1 ตัวอย่างของคำบรรยายที่ปรากฏในวีดิทัศน์ ก) คำบรรยายที่ปรากฏบริเวณด้านล่างของวีดิทัศน์ และ ข) คำบรรยายที่ปรากฏขึ้นในบริเวณที่น่าสนใจ
ที่มา: วีดิทัศน์ประกอบจาก Youtube ช่อง echo และ Mushroom TV

ทั้งนี้ โดยส่วนมากคำบรรยายที่ปรากฏบริเวณด้านล่างของวิดีโอจะเป็นแบบที่แสดงออกมาตามเสียงของผู้พูด หรือผู้บรรยายในวิดีโอ ทำให้สามารถใช้วิธีการทางด้านการตรวจจับตัวอักษร (Text Detection) เพื่อตรวจจับ บริเวณที่เป็นคำบรรยายที่ปรากฏในวิดีโอ และใช้วิธีการรู้จำเสียง (Speech Recognition) [1] ซึ่งเป็นงานวิจัยทาง ด้านปัญญาประดิษฐ์ (Artificial Intelligence) ดังนั้น ในวิทยานิพนธ์ฉบับนี้ ได้มุ่งเน้นในส่วนของการตรวจจับคำบรรยาย (Subtitle Detection) [2] เพื่อตรวจจับคำบรรยายทั้งแบบที่ปรากฏบริเวณด้านล่างของวิดีโอและปรากฏขึ้นในบริเวณที่น่าสนใจ (ดังภาพประกอบที่ 1) ซึ่งเป็นงานวิจัยที่เกี่ยวข้องกับรูปภาพ (Image) ผลลัพธ์ที่ได้จากการตรวจจับตัวอักษร แสดงดังภาพประกอบที่ 1.2ก นอกจากการตรวจจับคำบรรยายที่ปรากฏในวิดีโอ วิทยานิพนธ์ฉบับนี้ ได้ทำวิจัยในส่วนของการรู้จำคำบรรยาย (Subtitle Recognition) ซึ่งเป็นกระบวนการแปลงรูปภาพตัวอักษรให้เป็นข้อความที่อยู่ในรูปแบบของตัวอักษร (Text) ผลลัพธ์ที่ได้จากการรู้จำคำบรรยายแสดงดังภาพประกอบที่ 1.2

ถ้าไม่มีใจรัก อย่าไปเลี้ยงมัน มันเป็นการจับสัตว์มาทรมานในกรงใจ

สำเร็จ

ก) ผลลัพธ์: ถ้าไม่มีใจรัก อย่าไปเลี้ยงมัน มันเป็นการจับสัตว์มาทรมานในกรงใจ

ข) ผลลัพธ์: สำเร็จ

ภาพประกอบที่ 1.2 ตัวอย่างของการรู้จำคำบรรยาย ก) ภาพที่ได้จากขั้นตอนการตรวจจับคำบรรยาย และ ข) ผลลัพธ์ที่ได้จากการรู้จำคำบรรยาย โดยผลลัพธ์จะอยู่ในรูปแบบของข้อความ (Text)

พหุ ประถมศึกษา

จากที่ได้กล่าวมาข้างต้น ในวิทยานิพนธ์นี้มุ่งเน้นในการนำการเรียนรู้เชิงลึก (Deep Learning) มาเพื่อช่วยในการตรวจจับและรู้จำคำบรรยายวิดีโอ เนื่องจากเป็นวิธีการที่มีความถูกต้องและแม่นยำ ดังนั้น ในกระบวนการของการตรวจจับคำบรรยายวิดีโอ ได้ทดสอบกับ 3 วิธี ประกอบด้วย You only look once version 3 (YoloV3) [3], Tiny-YOLOv3 [3] และ RetinaNet [4] สำหรับการรู้จำคำบรรยายวิดีโอ ได้เลือกใช้วิธี Convolutional Neural Networks (CNN) ร่วมกับ Long short-Term Memory (LSTM) [5] และ Gated Recurrent Unit (GRU) [6] โดยข้อมูลวิดีโอที่ใช้ในวิทยานิพนธ์นี้ประกอบไปด้วยวิดีโอที่มีคำบรรยายเป็นภาษาไทย ภาษาอังกฤษ ตัวเลขไทย และตัวเลขอารบิก แสดงดังภาพประกอบที่ 1.3 ซึ่งได้เก็บรวบรวมจากวิดีโอทั้งหมดทั้งสิ้น 24 วิดีโอ โดยการวัดประสิทธิภาพในการตรวจจับคำบรรยายวิดีโอใช้ค่าอัตราความแม่นยำเฉลี่ย (Mean Average Precision: mAP) [7] และการรู้จำคำบรรยายวิดีโอ ใช้ค่าความผิดพลาดระดับตัวอักษร (Character Error Rate: CER) [8]



ภาพประกอบที่ 1.3 ตัวอย่างของคำบรรยายวิดีโอที่ใช้ในวิทยานิพนธ์

1.2 ความมุ่งหมายของการวิจัย

นำเสนอและเปรียบเทียบกระบวนการในการตรวจจับและรู้จำคำบรรยายวิดีโอ (Video Subtitle Detection and Recognition) ด้วยวิธีการเรียนรู้เชิงลึก (Deep Learning) เพื่อนำไปใช้ในการเก็บข้อมูลจากคำบรรยายในอนาคต

1.3 ขอบเขตของการวิจัย

1.3.1 ข้อมูลที่ใช้ในการวิจัย

1.3.1.1 รวบรวมวิดีโอทั้งหมดจำนวน 24 วิดีโอ

ที่มีคำบรรยายภาษาไทย ภาษาอังกฤษ ตัวเลขไทย และตัวเลขอารบิก วิดีโอที่ใช้ได้มาจาก Facebook [Thairath - ไทยรัฐออนไลน์\(3นาทิตั้ง\)](#), [TEP - Thailand Education Partnership ภาควิชาเพื่อการศึกษาไทย](#) Youtube [Bearhug](#), [KLUAYTHAI](#), [KRIT Eighth](#)

[Floor](#), [Genierock](#), [Taj Tracks](#), [Cakes & Eclairs](#), [7clouds](#), [DopeLyrics](#), [Equilanora](#), [JEMIN Apollo](#), [San Ko](#), [Tangerine JJY](#), [SNH48 Lyrics](#), [Lemoring](#), [Zaty Farhani](#)

1.3.1.2 เตรียมข้อมูลเพื่อนำไปใช้ในกระบวนการตรวจจับคำบรรยายวิดีโอ

โดยแปลงคลิปวิดีโอทั้งหมดเป็นรูปภาพ โดยรูปภาพที่ได้จากการแปลงมีขนาด 1280x720 พิกเซล (Pixel) เพื่อนำไปสร้าง Ground Truth ด้วยโปรแกรม labelImg ซึ่ง Ground Truth คือการกำหนดบริเวณที่เป็นคำบรรยาย (Subtitle) ที่ปรากฏในรูปภาพ (ดังภาพประกอบที่ 1.4) รูปภาพที่ใช้ในการทดลองประกอบด้วย 2,700 รูปภาพ



ภาพประกอบที่ 1.4 ตัวอย่างการสร้าง Ground Truth ในบริเวณที่เป็นคำบรรยายวิดีโอ

1.3.1.3 เตรียมข้อมูลเพื่อนำไปใช้ในกระบวนการรู้จำคำบรรยายวิดีโอ

โดยการสร้างคำตอบ (Label) ให้กับรูปภาพของคำบรรยาย ซึ่งคำตอบคือข้อความของคำบรรยายที่ปรากฏในวิดีโอ ดังภาพประกอบที่ 1.5 รูปภาพและคำบรรยายที่ใช้ในการทดลองประกอบด้วย 4,224 รูปภาพ

ตัวระเบิดเคลย์มอร์ หรือ ทุ่นระเบิดสังหารบุคคล

964_ตัวระเบิดเคลย์มอร์ หรือ ทุ่นระเบิดสังหารบุคคล

ภาพประกอบที่ 1.5 ตัวอย่างการเก็บข้อมูลคำบรรยายวิดีโอ

964_ตัวระเบิดเคลย์มอร์ หรือ ทุ่นระเบิดสังหารบุคคล [964 คือหมายเลขลำดับรูปภาพ และ

“ตัวระเบิดเคลย์มอร์ หรือ ทุ่นระเบิดสังหารบุคคล” คือคำตอบ (Label)]

1.3.3.2 ข้อมูลที่ใช้ในการรู้จำคำบรรยายวิดีโอ

รูปภาพคำบรรยายวิดีโอจำนวน 4,224 รูปภาพจะถูกแบ่งด้วยวิธี Cross-Validation และแบ่งเป็น 5-fold เพื่อใช้สำหรับตรวจสอบประสิทธิภาพในการรู้จำคำบรรยายวิดีโอ ด้วยค่าความผิดพลาดระดับตัวอักษร (Character Error Rate: CER)

1.3.4 กระบวนการเรียนรู้เชิงลึกสำหรับการตรวจจับและรู้จำคำบรรยายวิดีโอ

1.3.4.1 ตรวจจับคำบรรยายวิดีโอ ด้วยวิธีการเรียนรู้เชิงลึก

ได้แก่ YoloV3, Tiny-YOLOv3 และ RetinaNet

1.3.4.2 รู้จำคำบรรยายวิดีโอ ด้วยวิธีการเรียนรู้เชิงลึก

ได้แก่ CNN, LSTM และ GRU

1.4 ผลที่คาดว่าจะได้รับจากงานวิจัยครั้งนี้

กระบวนการเรียนรู้เชิงลึกสำหรับการตรวจจับและรู้จำคำบรรยายวิดีโอที่มีความถูกต้องสูงที่สามารถตรวจจับและรู้จำคำบรรยายที่มีขนาดและตำแหน่งแตกต่างกัน รวมไปถึงภาษาที่แตกต่างกัน

1.5 นิยามศัพท์เฉพาะ

การตรวจจับคำบรรยายวิดีโอ (Video Subtitle Detection) คือ การใช้โครงข่ายประสาทคอนโวลูชันในการตรวจจับตำแหน่งของคำบรรยาย โดยใช้วิธีต่าง ๆ เช่น YOLOv3, Tiny-YOLOv3 และ RetinaNet

รู้จำคำบรรยายวิดีโอ (Video Subtitle Recognition) คือ การใช้โครงข่ายประสาทแบบคอนโวลูชันในการพยากรณ์ข้อความของคำบรรยายโดยใช้วิธี เช่น VGG16 และ VGG19

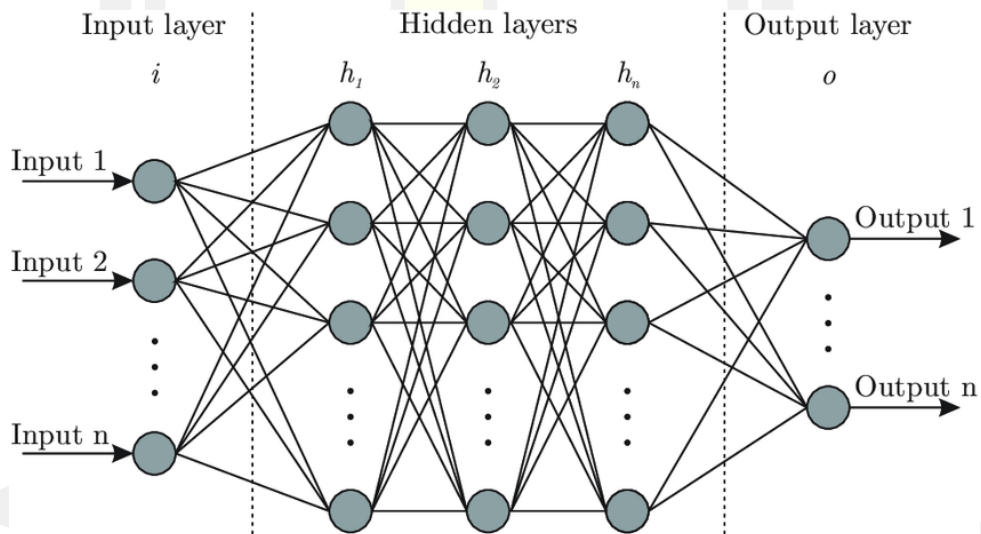
การเรียนรู้เชิงลึก (Deep Learning) คือ การเรียนรู้ของคอมพิวเตอร์ที่มีหลักการมาจากการทำงานของสมองมนุษย์ โดยจะเป็นการรับข้อมูลเข้ามา แล้วประมวลผลซ้ำ ๆ ในรูปแบบต่าง ๆ จำนวนมาก เพื่อทำการหาจุดที่จะบ่งบอกคุณลักษณะของข้อมูลที่นำเข้าไปได้ และแสดงผลลัพธ์ออกมาเป็นกลุ่มต่าง ๆ ว่าข้อมูลที่ใส่เข้าไปนั้นเป็นข้อมูลที่อยู่ในกลุ่มไหน เช่น การใส่รูป สุนัข เข้าไปเมื่อทำการประมวลผล แล้วผลลัพธ์จะออกมาเป็น สุนัข หรือ สัตว์ เป็นต้น

บทที่ 2

ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

2.1 การเรียนรู้เชิงลึก (Deep Learning)

การเรียนรู้เชิงลึกเป็นส่วนหนึ่งของการเรียนรู้ของเครื่องจักร (Machine Learning) [9] มีพื้นฐานมาจากโครงข่ายประสาทเทียม (Artificial Neural Network : ANN) เป็นวิธีที่สร้างขึ้นเพื่อให้เครื่องจักรสามารถเรียนรู้ได้โดยใช้ต้นแบบจากระบบประสาทของมนุษย์ โดยต้องใส่ข้อมูลเข้าไปในชั้นรับข้อมูล (Input Layer) จากนั้นเครื่องจักรจะนำข้อมูลไปประมวลผลในชั้นซ่อน (Hidden Layer) แล้วจะนำเสนอข้อมูลผลลัพธ์ในชั้นแสดงผล (Output Layer) ตามภาพประกอบที่ 2.1 โดยโครงข่ายประสาทเทียมได้มีการพัฒนาอย่างต่อเนื่องซึ่งวิธีที่นิยมใช้ในปัจจุบันคือ โครงข่ายประสาทคอนโวลูชัน (Convolution Neural Network: CNN) Long Short-Term Memory (LSTM) และ Gated Recurrent Unit (GRU)

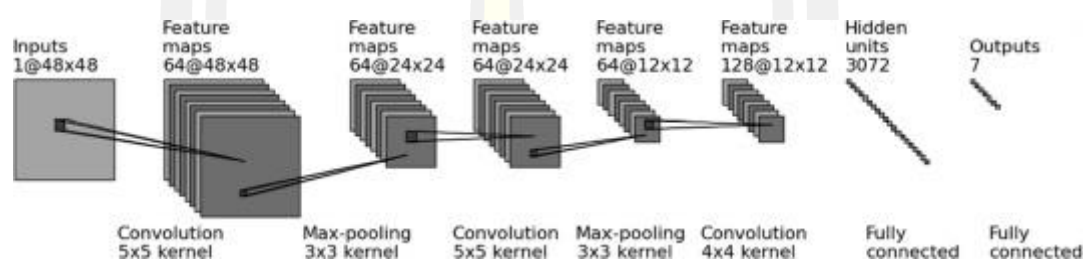


ภาพประกอบที่ 2.1 การเรียนรู้เชิงลึก (Deep Learning) [10]

2.2 โครงข่ายประสาทคอนโวลูชัน (Convolution Neural Network: CNN)

โครงข่ายประสาทคอนโวลูชัน (CNN) เป็นที่รู้จักในเรื่องของการเรียนรู้เชิงลึก (Deep Learning) ซึ่งมีการใช้ในงานวิจัยหลาย ๆ ด้าน เช่น การจัดกลุ่มข้อมูล การแบ่งรูปภาพ และการรู้จำคำพูด ในหลาย ๆ ปีมานี้มีสถาปัตยกรรม CNN จำนวนมากได้ถูกเสนอขึ้น เช่น EfficientNet [24], InceptionResNet [25] และ NASNet [26] และต่อมาในปี 2015 ได้มีการพัฒนา Visual

Geometry Group (VGG) โดย Simonyan และ Zisserman [27] จากมหาวิทยาลัยออกซฟอร์ด จุดประสงค์เพื่อตั้งชั้นของ CNN โดยเรียกว่า VGGNet ซึ่งได้ใช้ชั้นคอนโวลูชัน 16-19 ชั้น มาประมวลผลด้วยตัวกรองขนาดเล็กของชั้นคอนโวลูชัน โดยที่ขนาดของข้อมูลนำเข้า (input) มีขนาด 224x224 พิกเซล โดยที่แต่ละชั้นขนาดของข้อมูลจะถูกลดขนาดลงโดยใช้กระบวนการ max pooling โดยที่ขนาดของชั้นที่เล็กที่สุดมีขนาด 7x7 พิกเซล หลังจากนั้นจะตามมาด้วยชั้นเชื่อมโยงสมบูรณ์ (Fully Connected: FC) ซึ่งมีขนาด 4,096 4,096 1,000 ตามลำดับ และสิ้นสุดที่ softmax เป็นชั้นสุดท้ายตามภาพประกอบที่ 2.2



ภาพประกอบที่ 2.2 สถาปัตยกรรมโครงข่ายประสาทคอนโวลูชัน (Convolution Neural Network:

CNN) [11]

2.2.1 ชั้นรับข้อมูล (Input Layer)

ชั้นรับข้อมูลเป็นชั้นแรกของระบบการทำงานโดยจะทำหน้าที่ในการรับข้อมูลรูปภาพ รับขนาดความสูง (Height) ความกว้าง (Width) และความลึก (Depth) ของรูปภาพ โดยความสูงและความกว้างจะมีหน่วยเป็นพิกเซล และความลึกจะเป็นเลขตามสี เช่น ภาพสีเทา (Gray Scale) จะมีค่าเท่ากับ 1 และภาพสี (BGR) จะมีค่าเท่ากับ 3

พหุ ประถมศึกษา

2.2.2 ชั้นคอนโวลูชัน (Convolution Layer)

ชั้นคอนโวลูชันเป็นชั้นแรกถัดจากชั้นรับข้อมูล โดยจะทำหน้าที่ในการแยกคุณสมบัติ (Feature Extract) เช่น สี ขอบ รูปทรง จากข้อมูลที่ได้รับมาจากชั้นรับข้อมูล โดยจะทำการเปรียบเทียบรูปภาพที่รับจากชั้นรับข้อมูลกับตัวกรอง (Filter) โดยที่ขนาดของตัวกรองจะสามารถปรับได้ตามความเหมาะสมแต่จะนิยมใช้ที่ขนาด 3×3 สำหรับภาพที่มีขนาดไม่ใหญ่มาก 5×5 สำหรับภาพที่มีขนาดใหญ่ขึ้นมาโดยจะทำงานโดยการนำตัวเลขในขนาดของ filter มาคูณกับตัวเลขในขนาดของรูปภาพที่ตำแหน่งตรงกับ filter และทำการเลื่อนจากซ้ายไปขวา บนลงล่าง ตามภาพประกอบที่ 2.3

1	3	4	1	2	3
6	2	5	4	5	6
7	8	9	7	8	9
5	3	8	5	3	8
6	7	4	4	6	7
1	2	9	2	1	9

 \times

2	1	-2
0	0	-1
1	0	1

 $=$

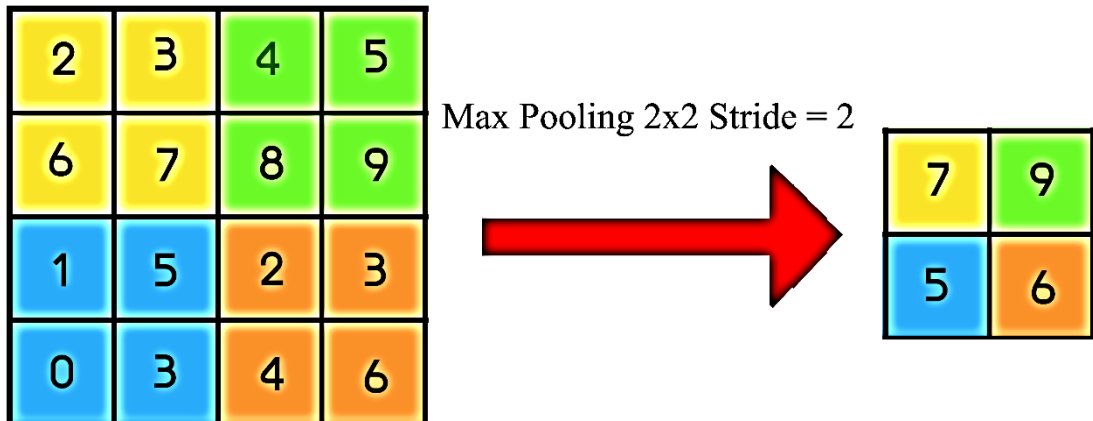
8			

$$1 \times 2 + 2 \times 1 + 3 \times (-1) + 4 \times 0 + 5 \times 0 + 6 \times (-1) + 7 \times 1 + 8 \times 0 + 9 \times 1 = 8$$

ภาพประกอบที่ 2.3 ตัวอย่างการทำงานของชั้นคอนโวลูชัน ขนาด 3×3 stride=1 [11]

2.2.3 ชั้นพูลลิ่ง (Pooling)

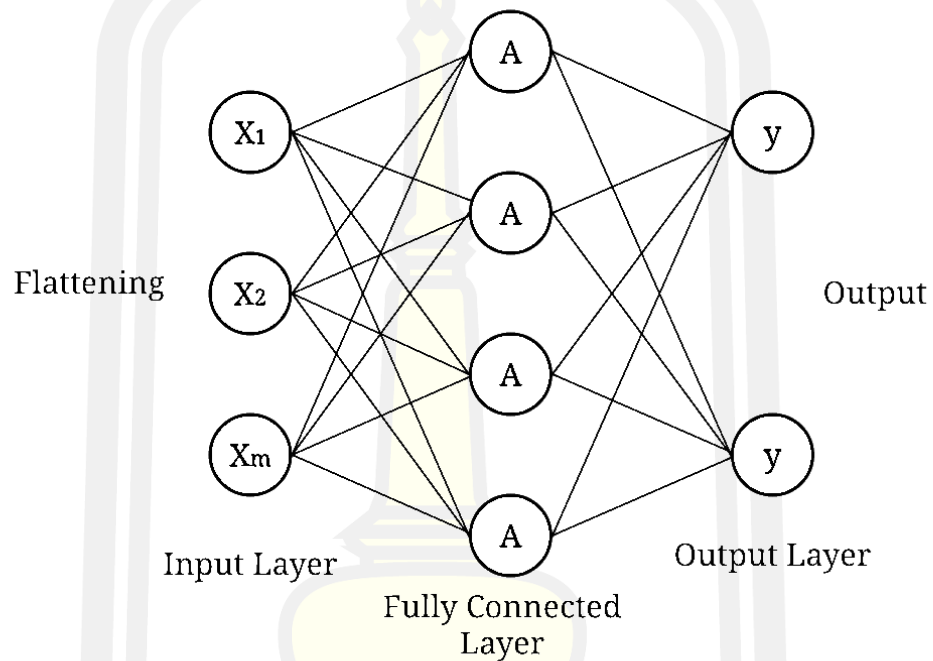
ชั้นพูลลิ่งเป็นชั้นที่จะต่อกับชั้นคอนโวลูชัน [11] ซึ่งจะทำหน้าที่ดึงค่าที่ต้องการคือค่าสูงสุด (Max Pooling) หรือ ค่าเฉลี่ย (Average Pooling) จากชั้นคอนโวลูชันเพื่อลดขนาดของข้อมูลลงตามขนาดดึงค่า (Pool Size) แต่ยังคงไว้ซึ่งลักษณะเด่นของข้อมูล



ภาพประกอบที่ 2.4 ตัวอย่างการทำงานของชั้น Max Pooling [11]

2.2.4 ชั้นเชื่อมโยงสมบูรณ์ (Fully Connected Layer)

ชั้นเชื่อมโยงสมบูรณ์เป็นชั้นที่เชื่อมต่อโดยสมบูรณ์กับชั้นต่าง ๆ ตามภาพประกอบที่ 2.5 ซึ่งจะมีลักษณะเป็น 1 มิติ แล้วนำไปสู่ชั้นผลลัพธ์ดังนั้นจึงเป็นชั้นสุดท้ายเพราะไม่สามารถนำไปเข้าชั้นคอนโวลูชัน หรือ ชั้นพูลลิงได้อีก



ภาพประกอบที่ 2.5 ตัวอย่างการเชื่อมต่อระหว่างชั้นรับข้อมูล ชั้นเชื่อมโยงสมบูรณ์ และชั้นผลลัพธ์ [12]

2.2.5 ชั้นผลลัพธ์ (output)

ชั้นผลลัพธ์เป็นชั้นที่ประมวลผลผลลัพธ์ที่ได้จากชั้นซ่อนออกมาแสดงผลเป็นตัวเลขที่ระบุตามกลุ่มของคำตอบ

2.3 การตรวจจับคำบรรยาย (Subtitle Detection)

การตรวจจับคำบรรยายคือการเรียนรู้ (Train) คำบรรยายจากตัวอย่างรูปภาพ และ นำมาผ่านขั้นตอนต่างๆเพื่อที่จะทำให้สามารถได้รับตำแหน่งคำบรรยายนั้นมา [13] โดยการตรวจจับในอดีตจะมุ่งเน้นไปที่การตรวจจับจากลักษณะของตัวอักษรซึ่งการตรวจจับได้มีการศึกษาและพัฒนาอย่างต่อเนื่องเพื่อที่จะตอบสนองต่อการตรวจจับสิ่งที่ต้องการจากรูปภาพนั้น ๆ ซึ่งเนื่องจากการตรวจจับคำบรรยายนั้นมีความซับซ้อนและแตกต่างกันมากไม่ว่าจะเป็น มุม ขนาด องศาและพื้นหลังที่ซับซ้อน [2] จึงได้มีการใช้โครงข่ายประสาทแบบคอนโวลูชัน (Convolutional Neural Network: CNN) เข้ามาช่วยในการตรวจจับลักษณะต่าง ๆ โดยใช้วิธีการ (Algorithm) เช่น You only look once Version 3 (YOLOv3), Tiny-YOLOv3 และ RetinaNet ซึ่งทั้ง 3 วิธีนี้จะถูกใช้ในส่วนของการเรียนรู้ตรวจจับตำแหน่งคำบรรยายภายในงานวิจัยนี้

YOLO เป็นสิ่งที่สร้างขึ้นมาเพื่อทำงานเกี่ยวกับการตรวจจับวัตถุรวมถึงการจัดกลุ่มวัตถุ โดยใช้วิธีการสร้างกรอบเพื่อกำหนดขอบเขตและจัดกลุ่มให้สิ่งที่อยู่ในขอบเขตนั้น โดย YOLO มีลักษณะคล้ายกับ Fully Convolutional neural network (FCNN) โดยจะทำการส่งรูปภาพผ่านไปสู่อ FCNN และได้รับผลลัพธ์ออกมาเป็นพิกัดที่ตรวจจับได้ โดย

2.3.1 YOLOv3

YOLOv3 เป็นวิธีที่พัฒนาต่อมาจาก YOLO โดยเป็นการเพิ่มความเร็วขึ้น 3 เท่า ซึ่ง YOLOv3 จะมีขั้นตอนการทำงานอยู่ 4 ขั้นตอน ขั้นแรกคือการคาดเดากรอบโดยจะคาดเดาออกมาเป็นตำแหน่งของ bounding box โดยมีค่าตำแหน่ง x, y ที่จะทำให้สามารถได้รับความกว้างและความสูงของตำแหน่งนั้น ตำแหน่งซึ่งจะสามารถวาดออกมาเป็นสี่เหลี่ยมเพื่อบ่งบอกตำแหน่ง ขั้นที่ 2 จะเป็นการคาดเดากลุ่มโดยใช้ binary cross-entropy แทนการใช้ softmax เนื่องจาก softmax มักจะเป็นการจำแนกประเภทได้เพียง 1 ประเภทต่อกรอบเท่านั้น แต่วิธีนี้จะสามารถจำแนกประเภทได้หลายประเภทเช่น ภาพที่มีผู้หญิง softmax มักจะเลือกว่าผลลัพธ์จะออกมาเป็นผู้หญิง หรือ มนุษย์ (Using a softmax imposes the assumption that each box has exactly one class which is often not the case) [3] เพียงอย่างเดียว แต่ binary cross-entropy จะจัดให้เป็นทั้งผู้หญิงและ มนุษย์ ขั้นที่ 3 เป็นการคาดเดาโดยใช้ขนาดที่แตกต่างกัน 3 ขนาด ขั้นที่ 4 เป็นการใช้ CNN 53 หรือเรียกอีกอย่างว่า Darknet-53 ซึ่งเป็นการใช้ CNN 53 ชั้น ตามภาพประกอบที่ 2.6 ชั้น ในการดึงคุณลักษณะเฉพาะ (Feature Extractor)

	Type	Filters	Size	Output
	Convolutional	32	3 x 3	256 x 256
	Convolutional	64	3 x 3 / 2	128 x 128
1 x	Convolutional	32	1 x 1	
	Convolutional	64	3 x 3	
	Residual			128 x 128
	Convolutional	128	3 x 3 / 2	64 x 64
2 x	Convolutional	64	1 x 1	
	Convolutional	128	3 x 3	
	Residual			64 x 64
	Convolutional	256	3 x 3 / 2	32 x 32
8 x	Convolutional	128	1 x 1	
	Convolutional	256	3 x 3	
	Residual			32 x 32
	Convolutional	512	3 x 3 / 2	16 x 16
8 x	Convolutional	256	1 x 1	
	Convolutional	512	3 x 3	
	Residual			16 x 16
	Convolutional	1024	3 x 3 / 2	8 x 8
4 x	Convolutional	512	1 x 1	
	Convolutional	1024	3 x 3	
	Residual			8 x 8
	Avgpool		Global	
	Connected		1000	
	Softmax			

ภาพประกอบที่ 2.6 สถาปัตยกรรม Darknet-53 [3]



2.3.2 Tiny-YOLOv3

Tiny-YOLOv3 เป็น YOLOv3 ที่ถูกตัดชั้นเชื่อมต่อสมบูรณ์ (Fully connected layer) และชั้นคอนโวลูชันออกไปบางส่วนโดยจะใช้เป็น Darknet19 ตามภาพประกอบที่ 2.7 หรือตัดออกเหลือเพียง 19 ชั้นทำให้ Tiny-YOLOv3 มีความเร็วที่มากกว่า YOLOv3 และมีความต้องการอุปกรณ์น้อยกว่า YOLOv3 ทำให้นิยมนำมาใช้ในระบบที่มีการตรวจจับตำแหน่งสิ่งต่าง ๆ แบบตามเวลาจริง (Real Time) แต่เนื่องจากการตัด ชั้นคอนโวลูชันออกไปบางส่วนทำให้มีประสิทธิภาพในการตรวจจับน้อยลงบ้าง

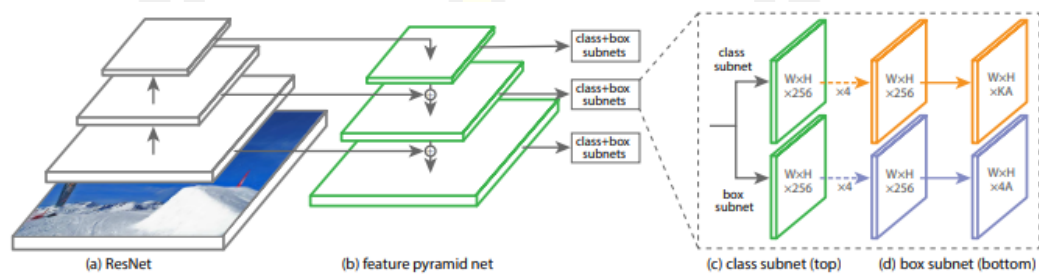
Layer	Type	Filters	Size/Stride	Input	Output
0	Convolutional	16	3 x 3/1	416 x 416 x 3	416 x 416 x 16
1	Maxpool		2 x 2/2	416 x 416 x 16	416 x 416 x 16
2	Convolutional	32	3 x 3/1	208 x 208 x 16	208 x 208 x 32
3	Maxpool		2 x 2/2	208 x 208 x 32	208 x 208 x 32
4	Convolutional	64	3 x 3/1	104 x 104 x 32	104 x 104 x 64
5	Maxpool		2 x 2/2	104 x 104 x 64	104 x 104 x 64
6	Convolutional	128	3 x 3/1	52 x 52 x 64	52 x 52 x 128
7	Maxpool		2 x 2/2	52 x 52 x 128	52 x 52 x 128
8	Convolutional	256	3 x 3/1	26 x 26 x 128	26 x 26 x 256
9	Maxpool		2 x 2/2	26 x 26 x 256	26 x 26 x 256
10	Convolutional	512	3 x 3/1	13 x 13 x 256	13 x 13 x 512
11	Maxpool		2 x 2/1	13 x 13 x 512	13 x 13 x 512
12	Convolutional	1024	3 x 3/1	13 x 13 x 512	13 x 13 x 1024
13	Convolutional	256	1 x 1/1	13 x 13 x 1024	13 x 13 x 256
14	Convolutional	512	3 x 3/1	13 x 13 x 256	13 x 13 x 512
15	Convolutional	255	1 x 1/1	13 x 13 x 512	13 x 13 x 255
16	YOLO				
17	Route 13				
18	Convolutional	128	1 x 1/1	13 x 13 x 256	13 x 13 x 128
19	Up-sampling			13 x 13 x 128	13 x 13 x 128
20	Route 19 8				
21	Convolutional	256	3 x 3/1	13 x 13 x 384	13 x 13 x 256
22	Convolutional	255	1 x 1/1	13 x 13 x 512	13 x 13 x 256
23	YOLO				

ภาพประกอบที่ 2.7 สถาปัตยกรรม Tiny-YOLOv3 [14]

พหุบัณฑิต ชีวะ

2.3.3 RetinaNet

RetinaNet เป็นการใช้ Feature Pyramid Network (FPN) [15] ซึ่งเป็นการตรวจจับข้อมูลที่จะปรับขนาดของภาพหลายอัตราส่วนเพื่อที่จะเพิ่มความแม่นยำในการตรวจจับแต่การทำเช่นนั้นทำให้อาจจะใช้เวลานานและต้องใช้หน่วยความจำ (memory) มากกว่าวิธีอื่น ตามภาพประกอบที่ 2.8 เป็นหลักในสถาปัตยกรรม เพื่อสร้างรูปภาพที่มีขนาดแตกต่างกัน โดย RetinaNet จะทำหน้าที่ 2 อย่างคือ การจัดกลุ่มให้กับตำแหน่งที่ล้อมกรอบว่าสิ่งต่าง ๆ ในกรอบคืออะไร และ วิเคราะห์กรอบที่ได้กับ ground truth ว่าตำแหน่งจุดพิกัดทั้ง 4 ตำแหน่งมีความสัมพันธ์กับ ground truth หรือไม่



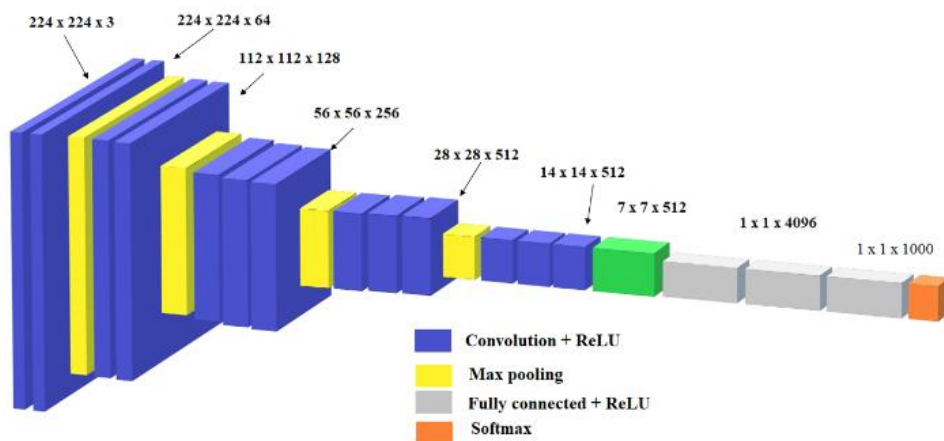
ภาพประกอบที่ 2.8 สถาปัตยกรรม RetinaNet [5]

2.4 การรู้จำคำบรรยาย (Subtitle Recognition)

การรู้จำคำบรรยายคือการเรียนรู้คำบรรยายจากตัวอย่างรูปภาพ และ นำมาผ่านขั้นตอนต่าง ๆ เพื่อที่จะทำให้สามารถได้รับคำบรรยายนั้นมา [13] ซึ่งคำบรรยายที่นำมาเพื่อเรียนรู้ในงานวิจัยนี้มีลักษณะที่เป็นประโยคของข้อความซึ่งจะมีความต่อเนื่องกันของคำบรรยายภายในประโยคดังนั้นจึงได้นำวิธี [16] Recurrent Neural Network (RNN) มาใช้ และนำมาถอดรหัส (Decode) โดยใช้ Connectionist Temporal Classification Loss (CTC Loss) [17]

2.4.1 Visual Geometry Group (VGG)

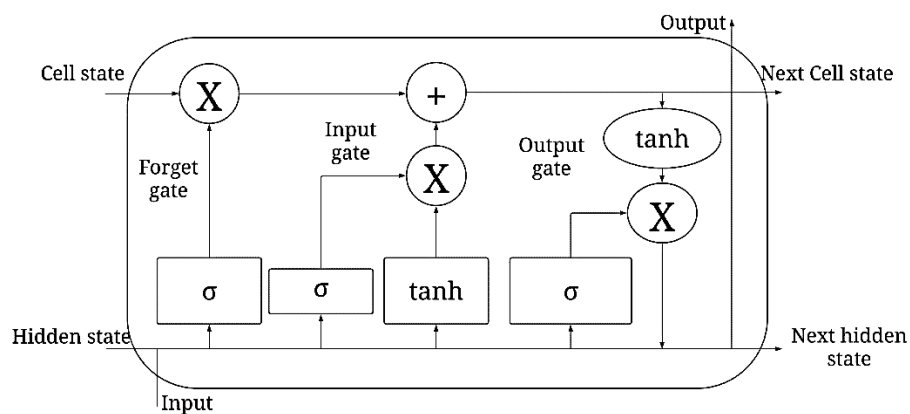
VGG [18] เป็นสถาปัตยกรรมที่นำ CNN ขนาด 3×3 stride=1 max pooling ขนาด 2×2 stride=2 fully connected layer 2 ชั้น และปิดด้วย softmax มาใช้ในการเรียนรู้เชิงลึกโดยที่将有ชั้นความลึก 16-19 ชั้น โดยส่วนมากจะนิยมใช้ที่ 16 ชั้นและ 19 ชั้น หรือเรียกว่า VGG16 และ VGG19 ตามลำดับ โดยที่สถาปัตยกรรม VGG16 สามารถดูได้ที่ ภาพประกอบที่ 2.9



ภาพประกอบที่ 2.9 สถาปัตยกรรม VGG16 [18]

2.4.2 Long Short-Term Memory (LSTM)

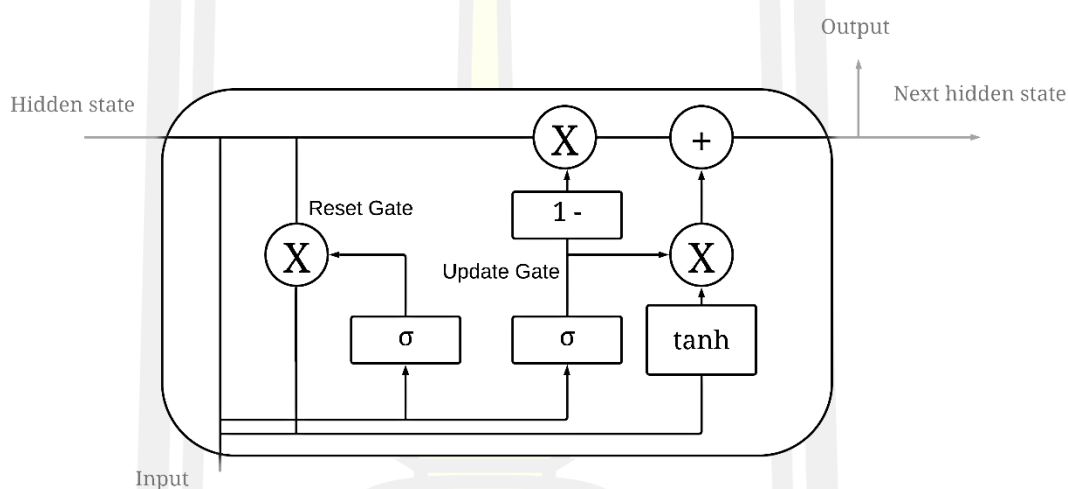
Hochreiter และ Schmidhuber [28] ได้นำเสนอถึง Recurrent Neural Network (RNN) ชนิดใหม่ซึ่งมีชื่อว่า Long Short-Term Memory (LSTM) ซึ่งได้นำมาช่วยในการแก้ไขปัญหา Vanishing Gradient ที่ถูกพบเมื่อใช้วิธี RNN กับข้อมูลที่มีความยาว เช่น เสียงพูดหรือวิดีโอ โดย LSTM ประกอบด้วย input gate, output gate และ forget gate ซึ่งจะเป็นสิ่งที่ควบคุมการไหลของข้อมูล โดยที่เมื่อ LSTM ได้รับข้อมูลมาจากชั้น input เป็นครั้งแรก LSTM จะเข้าสู่ input gate และเข้าสู่ output gate เพื่อตัดสินใจว่าจะเก็บค่าที่ได้ไว้แล้วจะวนซ้ำใน LSTM หรือแสดงผลข้อมูลหากเลือกที่จะวนซ้ำจะนำข้อมูลกลับมาเพื่อเข้าสู่ forget gate ซึ่งจะตัดสินใจว่าจะลบค่าที่เก็บไว้ทิ้งหรือยังคงเก็บค่าไว้ หากเก็บไว้จะไปรอการอัปเดตจาก input gate ซึ่งจะตัดสินใจว่าจะอัปเดตค่านั้นหรือไม่ และจะอัปเดตด้วยค่าอะไร แล้วส่งค่านั้นไป output gate เพื่อตัดสินใจว่าจะนำข้อมูลนั้นออกไปแสดงหรือจะนำข้อมูลกลับไปวนซ้ำอีกรอบ ดังนั้น LSTM จึงสามารถเรียนรู้จากข้อมูลที่เป็นลำดับและเก็บหรือลบข้อมูลทิ้งถ้าข้อมูลนั้นไม่จำเป็นตามภาพประกอบที่ 2.10



ภาพประกอบที่ 2.10 การทำงานของ LSTM

2.4.3 Gated Recurrent Unit (GRU)

GRU จะมีหลักการทำงานคล้ายกับ LSTM แต่มีความเร็วที่มากกว่าเพราะมีการตัด input และ output gate ออกในการวนซ้ำ เปลี่ยนมาใช้ reset gate และ update gate โดยที่ข้อมูล input เพื่อเข้ามาครั้งแรกจะทำการเก็บค่านั้นไว้แล้วจะตัดสินใจว่าจะนำข้อมูลไปวนซ้ำใน GRU หรือแสดงผล หากนำไปวนซ้ำจะเข้าสู่ reset gate เพื่อที่จะตัดสินใจว่าจะต้องลบข้อมูลไหน และเก็บข้อมูลไหนไว้ หลังจากนั้นจะเข้าสู่ update gate เพื่อที่จะตัดสินใจว่าจะอัปเดตข้อมูลไหน แล้วนำไปตัดสินใจว่าจะแสดงผลหรือนำไปวนซ้ำ ตามภาพประกอบที่ 2.11



ภาพประกอบที่ 2.11 การทำงานของ GRU

2.4.4 Connectionist Temporal Classification Loss (CTC Loss)

CTC Loss [17] เป็น loss ฟังก์ชันที่ถูกใช้ในการเรียนรู้ของ LSTM โดยเป็นการถอดรหัสเพื่อที่จะแก้ปัญหาของข้อมูลที่เป็นลำดับ ซึ่งจะสามารถคาดเดาข้อมูลที่ต่อเนื่องกันได้ ในช่วงนี้ได้ถูกใช้ในการจำแนกกลุ่มของข้อมูลที่เป็นลำดับ รวมไปถึงการจดจำลายมือและการจดจำเสียงพูด CTC จะค้นหาว่าข้อมูลนั้นเป็น blank หรือ ไม่มีตัวอักษร ดังนั้นข้อมูลที่เป็น blank หรือ ไม่มีตัวอักษรจะไม่ถูกเปลี่ยนไปเป็นตัวอักษรอื่น

2.5 การตรวจสอบประสิทธิภาพ

วิธีในการตรวจสอบความถูกต้องแม่นยำในการตรวจจับและรู้จำคำบรรยาย

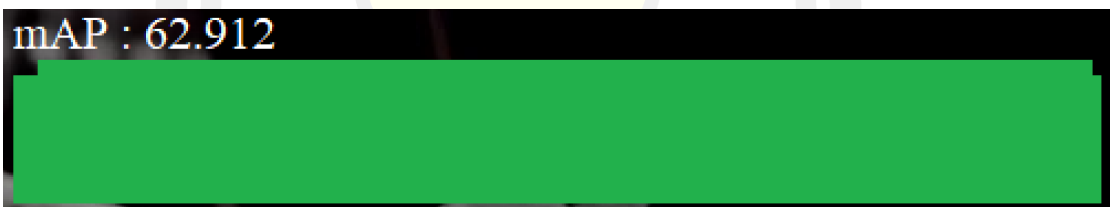
2.5.1 Intersection Over Union (IoU)

$$\frac{\text{Area of Overlap}}{\text{Area of Union}}$$

IoU เป็นวิธีที่ใช้ในการหาอัตราซ้อนทับกันของตำแหน่งที่ได้รับจากการตรวจจับตำแหน่งกับตำแหน่งที่กำหนดจากการทำ ground truth โดยการนำพื้นที่ที่ทับกัน (Area of Overlap) ตามภาพที่ 2.12 มาหารด้วยพื้นที่ทั้งหมด (Area of Union) ตามภาพที่ 2.13 โดยนิยมใช้ค่า IoU [19] ที่ 0.5 หรือข้อความที่ตรวจจับตรงกับตำแหน่งที่กำหนดใน ground truth มากกว่า 50% ตามภาพที่ 2.14



ภาพประกอบที่ 2.12 Area of Overlap โดยสีฟ้าคือพื้นที่ Ground Truth สีแดงคือผลลัพธ์ที่ได้จากการตรวจจับตำแหน่งคำบรรยาย สีเหลืองคือพื้นที่ซ้อนทับ (Area of Overlap)



ภาพประกอบที่ 2.13 Area of Union ซึ่งเป็นการรวมของพื้นที่ Ground Truth และผลลัพธ์การตรวจจับตำแหน่งคำบรรยาย



ภาพประกอบที่ 2.14 ค่า mAP 62.912 ซึ่งมากกว่ามีค่า IoU มากกว่าร้อยละ 50

2.5.2 อัตราความแม่นยำเฉลี่ย (mean Average Precision: mAP)

เป็นการหาค่าเฉลี่ยความถูกต้อง [7] ซึ่งได้จากการนำค่า Precision และ Recall มาทำเป็นกราฟและพื้นที่ใต้กราฟจะเป็นค่า Average Precision จากนั้นนำไปหาค่าเฉลี่ยจึงได้เป็นค่า mAP โดยที่ค่า Precision และ Recall หาได้ดังนี้

2.5.2.1 Precision

$$\frac{\text{TP}}{\text{TP} + \text{FP}}$$

เป็นค่าการคาดเดาที่ถูกต้องต่อการคาดเดาทั้งหมดโดยที่ tp คือค่าการคาดเดาที่ถูกต้องตาม ground truth fp คือการคาดเดาที่ผิดจาก ground truth

2.5.2.2 Recall

$$\frac{\text{TP}}{\text{TP} + \text{FN}}$$

เป็นค่าการคาดเดาที่ถูกต้องต่อการคาดเดาทั้งหมดโดยที่ tp คือค่าการคาดเดาที่ถูกต้องตาม ground truth fp คือการคาดเดาที่ผิดจาก ground truth

2.5.3 ความผิดพลาดระดับตัวอักษร (Character Error Rate: CER)

$$\frac{\text{Insertions} + \text{Substitutions} + \text{Deletions}}{\text{Length}}$$

จากสูตรการหาค่า CER [8] ด้านบนโดยที่ค่า Insertions คือจำนวนตัวอักษรที่นำเข้ามา Substitutions คือจำนวนตัวอักษรที่เปลี่ยนไป Deletions คือจำนวนตัวอักษรที่หายไป และ Length คือจำนวนตัวอักษรทั้งหมดที่ได้มาจากเฉลี่ยของ dataset

2.6 งานวิจัยที่เกี่ยวข้อง

2.6.1 งานวิจัยที่เกี่ยวข้องกับการตรวจจับคำบรรยาย (Text Detection)

Ma et al. [2] ได้สร้างวิธีตรวจจับข้อความแบบใหม่ชื่อว่า ReLaText ซึ่งเป็นเทคนิคใหม่ในการตรวจจับตำแหน่งคำบรรยายโดยใช้การเชื่อมความสัมพันธ์ (link relationship) ในการจัดกลุ่มข้อความ คอนโวลูชันแบบกราฟ (Graph Convolutional Network: GCN) ในการแยกความสัมพันธ์ที่ไม่ควรติดกันออก ซึ่ง ReLaText จะมีความแม่นยำในการตรวจจับข้อความที่มีช่องว่างหรือระยะห่างที่ใหญ่ (large inter-character) และแม่นยำในการตรวจจับข้อความที่มีระยะห่างระหว่างบรรทัดน้อย โดยใช้ RCTW-17, MSRA-TD500, Total-Text, CTW1500 และ DAST1500 ในการทดลอง ผลการทดลองด้วย RCTW-17 ด้วยรูปภาพขนาด 1500 พิกเซล พบว่าค่า Recall ได้ 61.7% F-Score ได้ 68.1% แต่มีค่า Precision 75.9% ซึ่งน้อยกว่าวิธี TextMountain อยู่ 4.9% เนื่องจากงานวิจัยนี้ได้อธิบายถึงการตรวจจับคำข้อความจึงได้นำมาอธิบายนั้นมาใช้เพื่อเป็นส่วนหนึ่งในการอธิบายถึงการตรวจจับคำบรรยาย

He et al. [20] ได้คิดวิธีที่จะลดเวลาในการตรวจจับข้อความที่มีหลายขนาด (multi-scale) โดยจะแบ่งออกเป็น 2 ขั้นตอน ซึ่งขั้นตอนแรกจะใช้ Scale-based Region Proposal Network (SRPN) ซึ่งมีประสิทธิภาพในการตรวจจับข้อความแบบกว้าง แล้วตัดพื้นที่ซึ่งไม่ใช่ข้อความออกรวมไปถึงปรับขนาดของข้อความ และเข้าสู่ขั้นตอนที่ 2 เป็นการใส่ เครือข่ายเชื่อมโยงสมบูรณ์ (Fully Convolutional Network) ซึ่งเป็นวิธีที่ใช้ ชั้นเพียงชั้นคอนโวลูชัน และ พูลลิง เท่านั้น โดยจะตัดชั้นเชื่อมโยงสมบูรณ์ออก (Fully connected layer) และใช้คอนโวลูชัน 1×1 แทนในการแปลงข้อความที่ตรวจจับมาจกขั้นตอนที่ 1 ซึ่งจะมีความแม่นยำสูงกว่าแต่่ามีความแคบในการทำงานน้อยกว่าขั้นตอนแรก เนื่องจากพื้นที่ซึ่งไม่ใช่ข้อความถูกนำออก และการปรับขนาดของข้อความทำให้วิธีนี้มีประสิทธิภาพและความเร็วที่สูง ซึ่งมีค่าวัดประสิทธิภาพ (F-measure) 85.40% ที่ 16.5 เฟรมต่อวินาที (frame per second: fps) และประสิทธิภาพในการแข่งขัน (competitive performance) 79.66% ที่ 35.1 fps เนื่องจากมีการอธิบายถึงความจำเป็นของการตรวจจับข้อความที่ต้องสามารถตรวจจับข้อความเดิมแบบที่มีความละเอียดสูงและแบบความละเอียดต่ำ ได้จึงได้หาวิธีที่สามารถตรวจจับแบบหลายขนาดได้

Wang et al. [21] ได้คิดวิธีในการเพิ่มประสิทธิภาพการตรวจจับตำแหน่งข้อความที่อยู่ห่างกันมากหรือที่มีอยู่ชิดกันมาก โดยการจับกลุ่มของตัวอักษรเพื่อเพิ่มประสิทธิภาพการตรวจจับข้อความที่มีลักษณะเป็นเส้นโค้งด้วยวิธี CNN โดยใช้สถาปัตยกรรม VGG16 เป็นหลักในการตรวจจับซึ่งรูปภาพที่นำมาใช้ในการเรียนรู้จะถูกทำ ground truth หลายรอบ โดยการหมุนข้อความตามเข็มนาฬิกาและทวนเข็มนาฬิกาการเพิ่มให้ข้อความแต่ละส่วนตั้งตรงจากการทดลองด้วย dataset DAST1500 ซึ่งเป็น dataset ที่เก็บรูปภาพที่มีลักษณะกระจัดกระจาย โค้ง หรือ แออัด และ dataset อื่น ๆ คือ ICDAR15, MTWL, Totaltext, SCUT-CTW1500 พบว่า 4 ใน 5 dataset วิธีการใหม่นี้มีประสิทธิภาพดีกว่าวิธีการเดิม แต่ยังมีข้อผิดพลาดในการตรวจจับข้อความที่ยากในการระบุ และข้อความที่มีขนาดเล็กเกินไป เนื่องจากที่มีการอธิบายถึงการทำให้ ground truth เพื่อกำหนดขนาดและตำแหน่งของข้อความตัวอย่างจึงได้มีการนำวิธี ground truth มาใช้

Zhu and Du [22] ได้สร้างวิธีใหม่ชื่อ TextMountain โดยจะเป็นการใช้ข้อมูลตำแหน่งของกรอบและจุดศูนย์กลางโดยการหาความน่าจะเป็นของจุดศูนย์กลางและกรอบโดยที่จุดศูนย์กลางจะเป็นเหมือนยอดภูเขาและส่วนบนและล่างของตัวอักษรจะเป็นดินเขา และหาทิศทางที่ข้อความจะไปเพื่อเพิ่มประสิทธิภาพในการตรวจจับ โดยได้มีการทดลองใช้ dataset ดังนี้ MLT, ICDAR2015, RCTW-17 และ SCUT-CTW1500 ซึ่งได้ใช้ ResNet-50 เป็นหลักในการทดลองในทุก dataset แต่จะพิเศษที่ ICDAR2015 ที่ได้มีการทดลองใช้ทั้ง ResNet-50 และ VGG-16 เนื่องจากได้มีการกล่าวถึง precision recall และ f-measure จึงได้ไปศึกษาและพบวิธีการแสดงผลการทดลองด้วยค่า mAP

Huang et al. [23] เป็นการเสนอวิธี Multiscale Connectionist Text Proposal Network (MS-CTPN) ที่จะสามารถตรวจจับข้อความแบบหลายขนาดและการเชื่อมโยงของข้อความ โดยใช้ Res-VGG16 เป็นพื้นฐาน ทดลองโดยใช้ข้อมูลจาก SynthText ซึ่งมีรูปภาพตัวอย่างทั้งหมด 12,000 รูป โดย 10,000 ใช้เรียนรู้ อีก 2,000 ใช้ทดสอบ ซึ่งผลการทดลองออกมาว่าการใช้ VGG16 + RPN จะใช้เวลาที่น้อยที่สุดคือ 0.2159 รูปต่อวินาที ค่า IoU 0.5624 หรือตำแหน่งที่ได้รับตรงกับตำแหน่ง ground truth 56.24% และ Res-VGG16 + Bi-RNN + MS-RPN ใช้เวลานานที่สุดแต่ค่า IoU 0.7635 หรือผลที่ได้ตรงกับ ground truth 76.35% ซึ่งเยอะที่สุด เนื่องจากได้อธิบายเกี่ยวกับ intersection over union (IoU) จึงได้นำไปศึกษาเพิ่มเติมและได้หาวิธีการกำหนด IoU เพื่อที่จะใช้ในการทดสอบ

2.6.2 งานวิจัยที่เกี่ยวข้องกับการรู้จำคำบรรยาย (Text Recognition)

Yan and Xu [24] ได้นำเสนอถึงสถาปัตยกรรม residual network โดยเป็นการใช้ Bi-GRU และ Connectionist temporal classification (CTC) ในการรู้จำคำบรรยายวิดิทัศน์ทั้งภาษาอังกฤษและภาษาจีน โดยสถาปัตยกรรมที่นำเสนอได้นั้นได้ทดลองกับ 2 dataset คือ ICDAR2003 และ ICDAR2013 ซึ่งมีค่าความแม่นยำ (Accuracy) 92.3% และ 89.2% ตามลำดับ

Jemni et al. [25] เป็นการตรวจจับและแก้ไขลายมือเขียนอารบิกที่ผิดหรือไม่มีในพจนานุกรม โดยการใช้ MDLSTM และ CNN ในการทดลองกับ dataset KHATT โดยได้แบ่งเป็น 2 ขั้นตอน ในขั้นแรก จะเป็นการตรวจจับตำแหน่งของข้อความ และในขั้นตอนที่สองจะเป็นการแก้ไขข้อความให้ถูกต้องโดยการใช้ข้อความและเลือกคำจากพจนานุกรม โดยใช้ต้นแบบ WSL, PSL, MSL เนื่องจากมีการทำทั้งตรวจจับข้อความและรู้จำข้อความจึงได้แบ่งการทดลองออกเป็นส่วนของ การตรวจจับคำบรรยายและการรู้จำคำบรรยาย

Gan et al. [26] เสนอถึง 1D-CNN และ Temporal Convolutional Recurrent Network (TCRN) โดยได้ตั้งชื่อว่า 1D-TCRN เพื่อที่จะนำมาใช้ในการรู้จำลายมือตัวอักษรจีนตาม IAHCCT dataset โดยนำมาเขียนในอากาศเพื่อที่จะรู้จำลำดับขีด และท่าทางการเขียน ในแบบจำลองนี้ นำ 1D residual convolution block มาใช้และนำมาเชื่อมกับ Sequence layer หรือก็คือเป็นการเชื่อมต่อของ CNN, LSTM และ CTC เพื่อนำมารู้จำลายมือตัวอักษรภาษาจีน 2,565 ตัวอักษร

Zhang et al. [27] ได้เสนอเทคนิคใหม่คือ Scale-aware hierarchical attention network for scene text recognition (SaHAN) ซึ่งได้รับแรงบันดาลใจจากเครือข่ายพีระมิด (pyramid network) โดยเป็นวิธีที่ใช้ทั้ง โครงข่ายประสาทคอนโวลูชัน (Convolution Neural Network: CNN) และ โครงข่ายประสาทแบบเกิดซ้ำ (Recurrent Neural Network : RNN) ซึ่งแบ่งเป็น 2 ส่วน ส่วนแรกคือการเข้ารหัสด้วย CNN และ RNN โดยจะใช้ ResNet-45 เป็นหลัก แล้วทำการลดขนาดลง (convolution layer) และอัดให้เหลือ 1 มิติ (fully connected layer) แล้ว

นำเข้าสู่ bidirectional long short-term memory (BiLSTM) และนำไปถอดรหัสด้วยการใช้ Gated Recurrent Unit (GRU) โดยวิธีนี้ได้ทดลองด้วย dataset 7 อย่าง คือ IIT5K-Words (IIT5K) , Street View Text (SVT) , ICDAR 2003 (IC03) , ICDAR 2013 (IC13) , ICDAR 2015 (IC15) , SVT-Perspective (SVT-P) , CUTE80 ได้ผลลัพธ์ดังนี้ IIT5K ที่ใช้พจนานุกรม 50 คำ มีความแม่นยำ 99.2% ที่ใช้พจนานุกรม 1000 คำ 97.9% ที่ไม่ใช้พจนานุกรม 91.2% SVT ที่ใช้พจนานุกรม 50 คำ 98.5% ที่ไม่ใช้พจนานุกรม 90.4% IC03 ที่ใช้พจนานุกรม 50 คำ 99.3% ทั้งหมด 98.4 ที่ไม่ใช้พจนานุกรม 95.5 IC13 ที่ไม่ใช้พจนานุกรม 93.0% IC15 ที่ไม่ใช้พจนานุกรม 75.0% IC15 แบบ 1811 รูป ที่ไม่ใช้พจนานุกรม 80.7% SVT-P ที่ใช้พจนานุกรม 50 คำ 95.2 ทั้งหมด 91.0 ที่ไม่ใช้พจนานุกรม 82.8 CUTE80 ที่ไม่ใช้พจนานุกรม 84.4% เนื่องจากการอธิบายถึง CNN, RNN, LSTM, BiLSTM และ GRU จึงช่วยลดระยะเวลาในการศึกษาข้อมูลเบื้องต้นทำให้เข้าใจข้อมูลเบื้องต้นเกี่ยวกับ CNN, RNN, LSTM, BiLSTM และ GRU

Chamchong et al. [28] เป็นการรู้จำลายมือเขียนไทยโบราณ โดยใช้ CNN, RNN และ CTC Loss มาสร้างสถาปัตยกรรม 6 แบบเพื่อเปรียบเทียบกัน โดยนำข้อมูลลายมือเขียนไทยโบราณมาจากห้องสมุดแห่งชาติไทย ซึ่งค่า CER ที่ดีที่สุดคือ 11.9% โดยจะใช้ตัวอักษร 44 ตัว สระ 18 ตัว วรรณยุกต์ 4 ตัว สัญลักษณ์อื่น ๆ 5 ตัว เลขไทย 10 ตัว รวมทั้งสิ้น 81 ตัว ในการทดสอบจะใช้ชุดข้อมูลที่เก็บไว้ก่อน ค.ศ. 1902 โดยมีทั้งหมด 140 รูป จากการทดสอบพบว่า การนำข้อมูลเข้าทดสอบพร้อมกัน 32 (batch 32) จะมีความผิดพลาดน้อยกว่าเข้าทีละ 64 (batch 64) LSTM จะมีความผิดพลาดมากกว่า GRU การมีชั้น CNN มากกว่าทำให้มีความผิดพลาดน้อยลง การใช้ CNN จะให้ผลลัพธ์ที่ดีกว่าการไม่ใช้ CNN เนื่องจากการทดลองในการรู้จำลายมือทั้งแบบใช้แต่ RNN รวมถึงใช้ทั้ง CNN และ RNN และได้ข้อสรุปว่า การใช้สถาปัตยกรรมที่มี CNN ให้ผลลัพธ์ที่ดีกว่าจึงได้สร้างสถาปัตยกรรมที่ใช้ CNN ร่วมกับ RNN

บทที่ 3

วิธีดำเนินการวิจัย

บทนี้ได้กล่าวถึง ขั้นตอนการดำเนินการวิจัยของการเรียนรู้เชิงลึกสำหรับการตรวจจับและรู้จำคำบรรยายวิดีโอ โดยกระบวนการทำวิจัยได้แบ่งออกเป็น 2 ขั้นตอนดังต่อไปนี้ การตรวจจับคำบรรยายวิดีโอ และการรู้จำคำบรรยายวิดีโอ ซึ่งได้อธิบายถึงรายละเอียดวิธีดำเนินการวิจัย ในหัวข้อดังต่อไปนี้

- 3.1 ชุดข้อมูลที่ใช้ในการวิจัย
- 3.2 การตรวจจับคำบรรยายวิดีโอและทดสอบประสิทธิภาพ
- 3.3 การรู้จำคำบรรยายวิดีโอและทดสอบประสิทธิภาพ

3.1 ชุดข้อมูลที่ใช้ในงานวิจัย

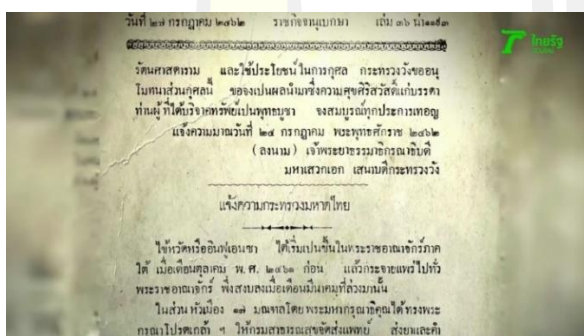
ข้อมูลวิดีโอที่ใช้ในวิทยานิพนธ์คือวิดีโอจำนวนทั้งสิ้น 24 วิดีโอ ที่เก็บรวบรวมจาก Youtube และ Facebook โดยคำบรรยายที่ปรากฏในวิดีโอประกอบไปด้วย ภาษาไทย ภาษาอังกฤษ ตัวเลขไทย และตัวเลข อารบิก หลังจากเก็บรวบรวมวิดีโอ จึงได้นำวิดีโอทั้งหมดไปแยกเป็นเฟรม (Frame) เพื่อนำไปทำ Ground Truth สำหรับใช้ในการทดสอบและวัดประสิทธิภาพ โดยเรียกชุดข้อมูลนี้ว่า ชุดข้อมูลรูปภาพคำบรรยาย (Video Subtitle Image Dataset) โดยชุดข้อมูลแบ่งออกเป็น 2 ส่วน ดังนี้

3.1.1 รูปภาพที่ใช้สำหรับตรวจจับคำบรรยายวิดีโอทัศน์

รูปภาพคำบรรยาย คือรูปภาพที่มีคำบรรยายปรากฏอยู่ 2 รูปแบบคือ คำบรรยายที่ปรากฏบริเวณด้านล่างของวิดีโอทัศน์ และคำบรรยายที่ปรากฏขึ้นในบริเวณที่น่าสนใจ ซึ่งคำบรรยายนั้นมีทั้งภาษาไทย ภาษาอังกฤษ ตัวเลขไทย และตัวเลขอารบิก รวมทั้งสิ้น 2,700 รูปภาพ ตัวอย่างของรูปภาพคำบรรยายแสดงดังภาพประกอบที่ 3.1



ก)



ข)



ค)

ภาพประกอบที่ 3.1 ตัวอย่างรูปภาพที่ใช้สำหรับตรวจจับคำบรรยายวิดีโอทัศน์ ก) รูปภาพที่ประกอบด้วยภาษาอังกฤษและเลขอารบิก ข) รูปภาพที่ประกอบด้วยภาษาไทยและเลขไทย ค) รูปภาพที่ประกอบด้วยภาษาไทยและเลขอารบิก

3.1.2 รูปภาพคำบรรยายวิดีโอที่ขึ้นที่ใช้สำหรับการรู้จำ

รูปภาพคำบรรยายวิดีโอที่ขึ้นที่ใช้สำหรับการรู้จำ เป็นรูปภาพคำบรรยาย (Subtitle Image) ที่นำมาจาก Ground Truth โดยนำมาเพื่อใช้สำหรับการรู้จำคำบรรยายวิดีโอ (Video Subtitle Recognition) โดยเป็นรูปภาพคำบรรยายที่เป็น ภาษาไทย ภาษาอังกฤษ ตัวเลขไทย และ ตัวเลขอารบิก จำนวนทั้งสิ้น 4,224 รูปภาพ ตัวอย่างของรูปภาพคำบรรยายวิดีโอแสดงดัง ภาพประกอบที่ 3.3ก โดยรูปภาพคำบรรยายทุกรูปจะถูกนำไปสร้างคำตอบ (Label) เพื่อใช้สำหรับ ตรวจสอบความถูกต้องในการรู้จำ และวัดประสิทธิภาพของอัลกอริทึม ตัวอย่างการสร้างคำตอบให้กับ รูปภาพคำบรรยาย วิดีโอที่ขึ้นแสดงดังภาพประกอบที่ 3.3ข

ส่วนคำกล่าวขานที่ว่า

1147_ ส่วนคำกล่าวขานที่ว่า

“เมื่อได้ลิ้มลองเนื้อมนุษย์แล้ว จะหากินอยู่เรื่อยๆ”

1148_ เมื่อได้ลิ้มลองเนื้อมนุษย์แล้ว จะหากินอยู่เรื่อยๆ

ก)

1147_ ส่วนคำกล่าวขานที่ว่า

1148_ เมื่อได้ลิ้มลองเนื้อมนุษย์แล้ว จะหากินอยู่

1147_ คือหมายเลขลำดับรูปภาพ

เรื่อยๆ

“ส่วนคำกล่าวขานที่ว่า” คือคำตอบ (Label)

1148_ คือหมายเลขลำดับรูปภาพ

“เมื่อได้ลิ้มลองเนื้อมนุษย์แล้ว จะหากินอยู่
เรื่อยๆ” คือคำตอบ (Label)

ข)

ภาพประกอบที่ 3.3 ตัวอย่างรูปภาพคำบรรยายวิดีโอที่ขึ้น ก) รูปภาพคำบรรยายวิดีโอ และ ข)

คำตอบ (Label)

พหุ ประถมศึกษา

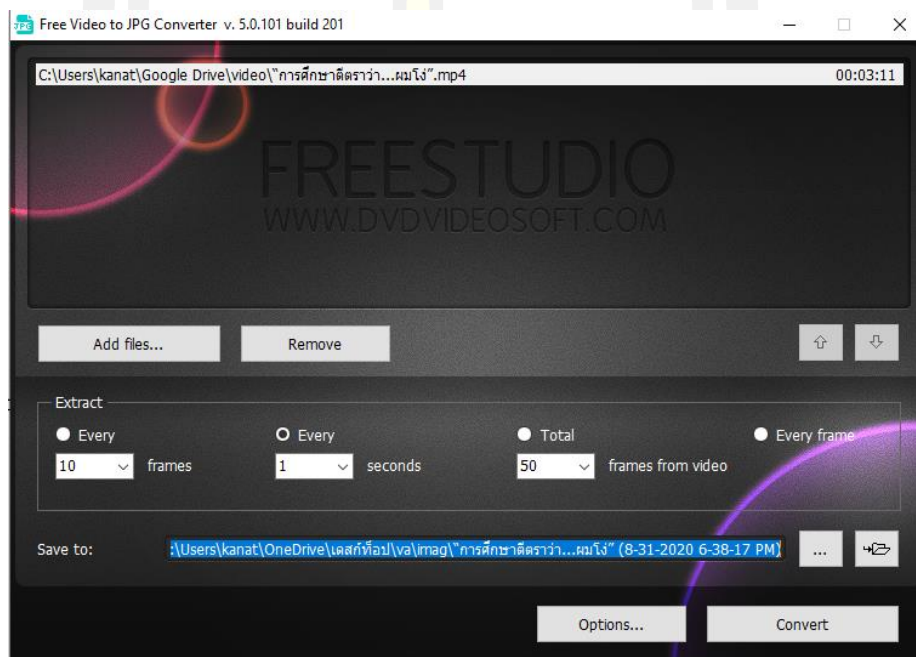
3.2 การเตรียมข้อมูล

การเตรียมข้อมูลแบ่งออกเป็น

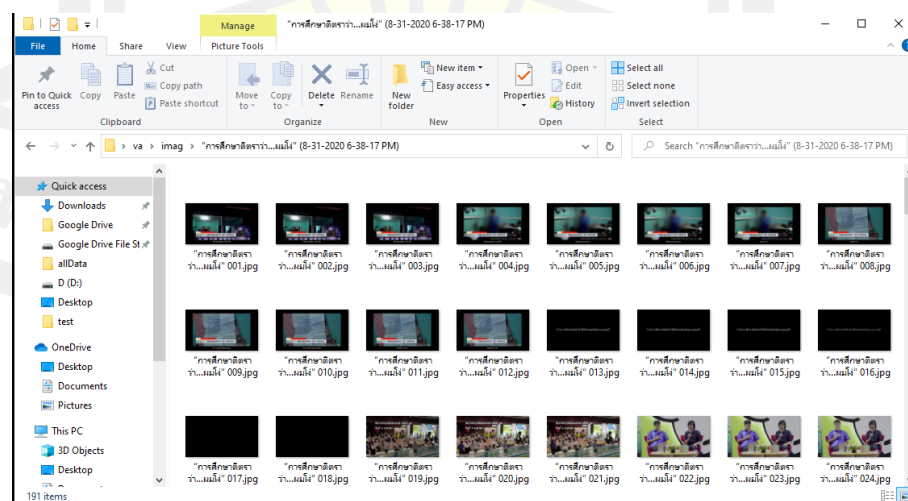
3.2.1 หาตัวอย่างที่มีการฝังคำบรรยายภาษาไทยและอังกฤษ

3.2.2 นำคลิปวิดีโอที่ค้นไปแปลงเป็นรูปภาพ

ด้วยโปรแกรม [Free Video to JPG Converter](#) ภาพประกอบที่ 3.4 ซึ่งเป็นโปรแกรมแปลงวิดีโอที่ค้นเป็นไฟล์รูปภาพ และจะได้ผลลัพธ์ออกมาเป็นรูปภาพ ดังภาพประกอบที่ 3.5



ภาพประกอบที่ 3.4 โปรแกรม Free Video to JPG Converter



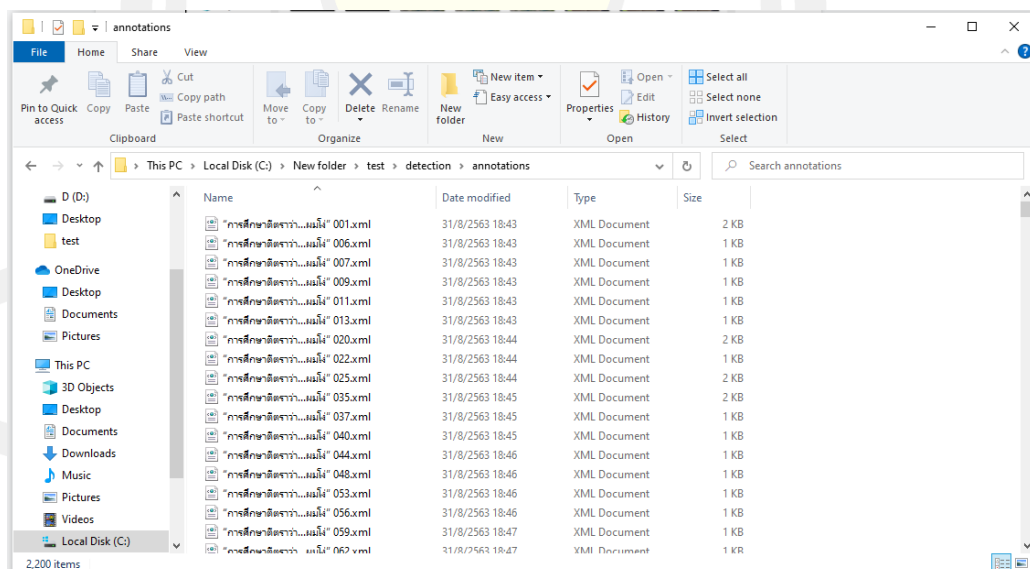
ภาพประกอบที่ 3.5 ภาพที่ได้จากโปรแกรม Free Video to JPG Converter

3.2.3 นำรูปภาพไปสร้าง ground truth

ด้วยโปรแกรม [Labellmg](#) ภาพประกอบที่ 3.6 และได้ผลลัพธ์ออกมาเป็นไฟล์ .xml ตามภาพประกอบที่ 3.7



ภาพประกอบที่ 3.6 ภาพโปรแกรม Labellmg



ภาพประกอบที่ 3.7 ภาพไฟล์ .xml ที่ได้จากโปรแกรม Labellmg

3.2.4 นำรูปภาพคำบรรยาย (Subtitle image) มาสร้างไฟล์ข้อความตัวอักษรประกอบ (Text transcription)

โดยการตัดมาเฉพาะพื้นที่ที่ถูกทำ ground truth และตั้งชื่อไฟล์ตามคำบรรยายในรูปภาพ ดังภาพประกอบที่ 3.8



ภาพประกอบที่ 3.8 ภาพไฟล์ข้อความคำบรรยาย

3.3 การตรวจจับคำบรรยายวิดีโอ

วิทยานิพนธ์นี้ได้ทำการทดลองการตรวจจับและการรู้จำคำบรรยายวิดีโอบน Google Colab ซึ่งประมวลผลด้วยหน่วยประมวลผลภาพ (Graphic Processing Unit: GPU) โดยกระบวนการตรวจจับคำบรรยายวิดีโอ แบ่งเป็นขั้นตอนดังนี้

3.3.1 ข้อมูลที่ใช้ในการทดสอบกระบวนการตรวจจับคำบรรยายวิดีโอ

ข้อมูลที่ใช้ในการทดสอบกระบวนการตรวจจับคำบรรยายวิดีโอคือ รูปภาพที่ใช้สำหรับตรวจจับคำบรรยายวิดีโอ จำนวนทั้งสิ้น 2,700 รูปภาพ โดยแบ่งข้อมูลที่ใช้ในการทดสอบด้วยวิธี k fold Cross-Validation โดยกำหนดให้ $k=5$

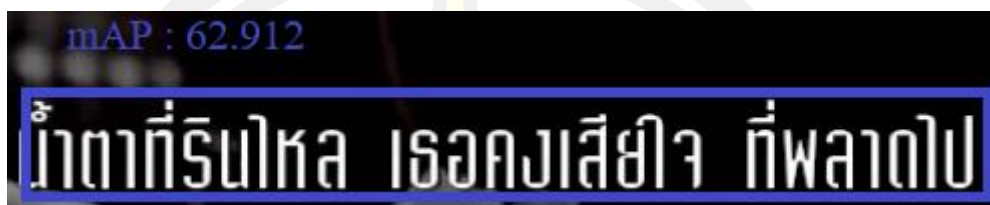
3.3.2 ทำการสร้างโมเดลในการตรวจจับคำบรรยายวิดีโอ

ทำการสร้างโมเดลในการตรวจจับคำบรรยายวิดีโอด้วยวิธีการเรียนรู้เชิงลึกด้วยวิธี YoloV3, Tiny-YOLOv3 และ RetinaNet

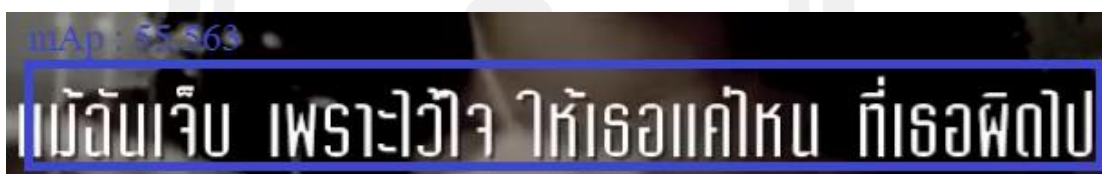
พหุบัณฑิต ชีวะ

3.3.3 ทดสอบประสิทธิภาพของโมเดลในการตรวจจับคำบรรยายวิดีโอทัศน์

ทดสอบประสิทธิภาพของโมเดลในการตรวจจับคำบรรยายวิดีโอทัศน์ด้วยค่าความแม่นยำเฉลี่ย (Mean Average Precision: mAP) โดยกำหนดให้มีค่า IoU ที่ 0.5 ตามภาพประกอบที่ 3.9 และ 3.10



ภาพประกอบที่ 3.9 ตัวอย่างผลลัพธ์ค่า mAP 62.912%



ภาพประกอบที่ 3.10 ตัวอย่างผลลัพธ์ค่า mAP 55.563%

3.4 การรู้จำคำบรรยายวิดีโอทัศน์

การรู้จำคำบรรยายวิดีโอทัศน์ แบ่งเป็นขั้นตอนดังนี้

3.4.1 ข้อมูลที่ใช้ในการทดสอบในการรู้จำคำบรรยายวิดีโอทัศน์

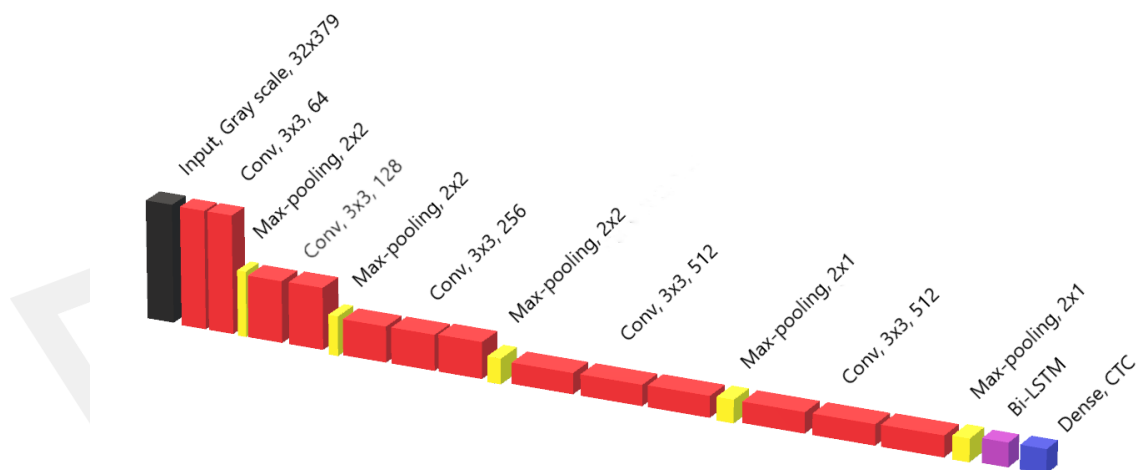
ข้อมูลที่ใช้ในการทดสอบในการรู้จำคำบรรยายวิดีโอทัศน์คือ รูปภาพคำบรรยายวิดีโอทัศน์ที่ใช้สำหรับการรู้จำ จำนวนทั้งสิ้น 4,224 รูปภาพ โดยแบ่งข้อมูลที่ใช้ในการทดสอบด้วยวิธี Cross-Validation โดยกำหนดให้ $k=5$

3.4.2 ทำการสร้างโมเดลในการรู้จำคำบรรยายวิดีโอทัศน์

ด้วยวิธีการเรียนรู้เชิงลึกด้วยวิธี Convolutional Neural Network (CNN), Long Short-Term Memory (LSTM) และ Gated Recurrent Unit (GRU) ตามตารางที่ 3.1, ตารางที่ 3.2, ตารางที่ 3.3 และภาพประกอบที่ 3.11

ตารางที่ 3.1 ตัวอย่างสถาปัตยกรรม CNN-LSTM

3.4.3 Stage	Operations	Resolution	Channels	Layers
1	Conv 3x3	32x379	64	2
2	Max-pooling 2x2			
3	Conv 3x3	16x189	128	2
4	Max-pooling 2x2			
5	Conv 3x3	8x94	256	3
6	Max-pooling 2x2			
7	Conv 3x3	4x94	512	3
8	Max-pooling 2x1			
9	Conv 3x3	2x94	512	3
10	Max-pooling 2x1			
11	Bi-LSTM	94	256	2
12	Dense & Softmax Function	158		1
13	CTC Loss Function			



ภาพประกอบที่ 3.11 สถาปัตยกรรม CNN-LSTM

ตารางที่ 3.2 สถาปัตยกรรมที่ใช้ในการสร้างแบบจำลองการรู้จำคำบรรยายวิดีโอ

Model A	Model B	Model C	Model D
16 weight layers	16 weight layers	19 weight layers	19 weight layers
Input (gray image), 32x379			
3x3 conv,64	3x3 conv,64	3x3 conv,64	3x3 conv,64
3x3 conv,64	3x3 conv,64	3x3 conv,64	3x3 conv,64
Max-pooling 2x2, stride 2			
3x3 conv,128	3x3 conv,128	3x3 conv,128	3x3 conv,128
3x3 conv,128	3x3 conv,128	3x3 conv,128	3x3 conv,128
Max-pooling 2x2, stride 2			
3x3 conv,256	3x3 conv,256	3x3 conv,256	3x3 conv,256
3x3 conv,256	3x3 conv,256	3x3 conv,256	3x3 conv,256
3x3 conv,256	3x3 conv,256	3x3 conv,256	3x3 conv,256
		3x3 conv,256	3x3 conv,256
Max-pooling 2x1, stride 1			
3x3 conv,512	3x3 conv,512	3x3 conv,512	3x3 conv,512
3x3 conv,512	3x3 conv,512	3x3 conv,512	3x3 conv,512
3x3 conv,512	3x3 conv,512	3x3 conv,512	3x3 conv,512
		3x3 conv,512	3x3 conv,512
Max-pooling 2x1, stride 1			
3x3 conv,512	3x3 conv,512	3x3 conv,512	3x3 conv,512
3x3 conv,512	3x3 conv,512	3x3 conv,512	3x3 conv,512
3x3 conv,512	3x3 conv,512	3x3 conv,512	3x3 conv,512
		3x3 conv,512	3x3 conv,512
BiLSTM	BiGRU	BiLSTM	BiGRU
BiLSTM	BiGRU	BiLSTM	BiGRU
Dense & Softmax Function,158			
CTC Loss Function			

ตารางที่ 3.3 สถาปัตยกรรมตามบทความของ Chamchong et al.

Model 1	Model 2
6 weight layers	6 weight layers
Input (gray image), 32x379	
Conv 3x3, 16	Conv 3x3, 16
Max-pooling 2x2	
Conv 3x3, 32	Conv 3x3, 32
Max-pooling 2x2	
Conv 3x3, 32	Conv 3x3, 32
Bi-GRU	Bi-LSTM
Bi-GRU	Bi-LSTM
Dense & Softmax Function, 157	

ตัวอย่างการทำงานของ Model A โดยเริ่มจากการนำข้อมูลขนาดความสูง 32 พิกเซล และความกว้าง 379 พิกเซลเข้าไปประมวลผลด้วยชั้น คอนโวลูชันชั้นแรกซึ่งมีขนาด 3x3 และมีค่า filter 64 และนำไปลดขนาดครึ่งหนึ่งโดยการตั้งค่าสูงสุดที่ได้จากชั้นคอนโวลูชัน แล้วจึงนำไปเข้าสู่ชั้นคอนโวลูชันและชั้นพลูลิงไปจนได้ขนาด none,93,512 โดยที่เลข 93 หมายถึงความยาวของลำดับ (sequence length) และ 512 คือ filter แล้วนำไปเข้าสู่ชั้น LSTM หรือ GRU ที่มี 2 ชั้น โดย LSTM หรือ GRU จะทำงานเป็นคู่ซึ่งแต่ละคู่จะเรียกว่า BiLSTM หรือ BiGRU ซึ่ง LSTM ชั้นแรกจะทำงานจากซ้ายไปขวา และ LSTM ชั้นที่ 2 จะทำงานจากขวาไปซ้ายแล้วจึงนำผลลัพธ์ที่ได้ไปรวมกันที่ชั้นถัดไปซึ่งก็คือชั้น Dense โดยในชั้นนี้ข้อมูลจะมีลักษณะดังนี้ None,93,158 ซึ่งเลข 93 คือ เลขของความยาวลำดับ (sequence length) และ 158 คือเลขของข้อมูลที่เป็นไปได้ทั้งหมดบวกหนึ่งโดยจะเป็นตัวอักษรภาษาอังกฤษตัวเล็ก ตัวอักษรภาษาอังกฤษตัวใหญ่ พยัญชนะไทย สระ วรรณยุกต์ และเลขไทย เพื่อใช้ในการหาคำตอบและนำไปถอดรหัสโดยใช้ CTC Loss ซึ่งเป็นหนึ่งในชั้น output โดยจะแปลงตัวเลขที่ได้จากชั้น Dense ให้เป็นตัวอักษรตามที่ได้กำหนดไว้

3.4.4 ทดสอบประสิทธิภาพของโมเดลในการรู้จำคำบรรยายวิดีโอ

ด้วยค่าความผิดพลาดระดับตัวอักษร (Character Error Rate: CER) โดยประโยคที่มีตัวอักษรเยอะที่สุดมี 90 ตัวอักษร ตัวอย่างผลลัพธ์ที่ได้จากการทดลอง ตามภาพประกอบที่ 3.12, ภาพประกอบที่ 3.13 และ ภาพประกอบที่ 3.14

ส่วนคำกล่าวขานที่ว่า

Original text = ส่วนคำกล่าวขานที่ว่า

Predicted text = ส่วนคำกล่าวขานที่ว่า

ภาพประกอบที่ 3.12 ตัวอย่างผลลัพธ์มีผลลัพธ์ตรงกับข้อมูลนำเข้าทั้งหมด

ดังนั้นจากสูตรการหาค่า $CER = \frac{I+S+D}{N}$ ดังนั้น $CER = \frac{0+0+0}{20} = 0.0$ หรือ
ผิดพลาด 0%

หลุมหลบภัยทางอากาศของทหารญี่ปุ่นเท่านั้น

Original text = หลุมหลบภัยทางอากาศของทหารญี่ปุ่นเท่านั้น

Predicted text = หลุมหลบภัยทางอากาศของทหารญี่ปุ่นเท่านั้น

ภาพประกอบที่ 3.13 ตัวอย่างผลลัพธ์โดยมีการเปลี่ยนจาก ๔ เป็น ๑ ตัวอักษร จากทั้งหมด 40
ตัวอักษร

ดังนั้นจากสูตรการหาค่า $CER = \frac{I+S+D}{N}$ ดังนั้น $CER = \frac{0+1+0}{40} = 0.025$
หรือผิดพลาด 2.5%

นำข้อมูลอันเป็นเท็จเข้าสู่ระบบคอมพิวเตอร์

Original text = นำข้อมูลอันเป็นเท็จเข้าสู่ระบบคอมพิวเตอร์

Predicted text = นำข้อมูลอันเป็นเท็จเข้าสู่ระบบคอมพิวเตอร์

ภาพประกอบที่ 3.14 ตัวอย่างผลลัพธ์โดยมีการเปลี่ยนจาก ท เป็น ก ที่ตำแหน่ง 16^๕ หายไปตรง
ตำแหน่งที่ 17 ทำให้หลังตำแหน่งที่ 17 จะเปลี่ยนไปทั้งหมดรวมถึง ' เป็น'

ดังนั้นจากสูตรการหาค่า $CER = \frac{I+S+D}{N}$ ดังนั้น $CER = \frac{1+23+0}{41} = 0.59$ หรือ
ผิดพลาด 59%

บทที่ 4

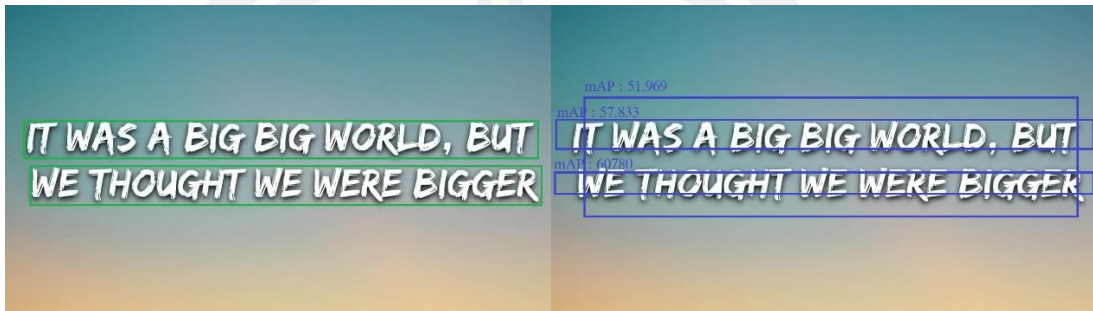
ผลลัพธ์การทดลอง

ในบทนี้แสดงถึงผลลัพธ์ที่ได้จากการทดลองโดยการใช้ TensorFlow framework ใน Google Colab และใช้ฟังก์ชันของ Google Colab ในการดูว่า GPU ที่ใช้เป็น GPU อะไร โดยเลือกใช้ GPU Tesla T4 ในการทดลองตรวจจับคำบรรยาย และ GPU Tesla P100 ในการรู้จำคำบรรยาย โดยใช้ dataset ของคำบรรยายวิดีโอที่ได้อาจจากการใช้วิธี 5-fold cross-validation โดยมีค่า mean Average Precision (mAP) ที่ Intersect over Union (IoU) = 0.5 ในการวัดประสิทธิภาพการตรวจจับคำบรรยายวิดีโอ ซึ่งหากค่า mAP มีค่ามากแสดงว่าการทดสอบนั้นมีผลลัพธ์ที่ดีกว่าแบบที่มีผลลัพธ์ที่มีค่าน้อยกว่า และค่า Character Error Rate (CER) เป็นค่าสำหรับวัดประสิทธิภาพของการรู้จำคำบรรยายวิดีโอ ถ้าหากค่า CER มีค่าน้อยแสดงว่าผลลัพธ์นั้นดีกว่า

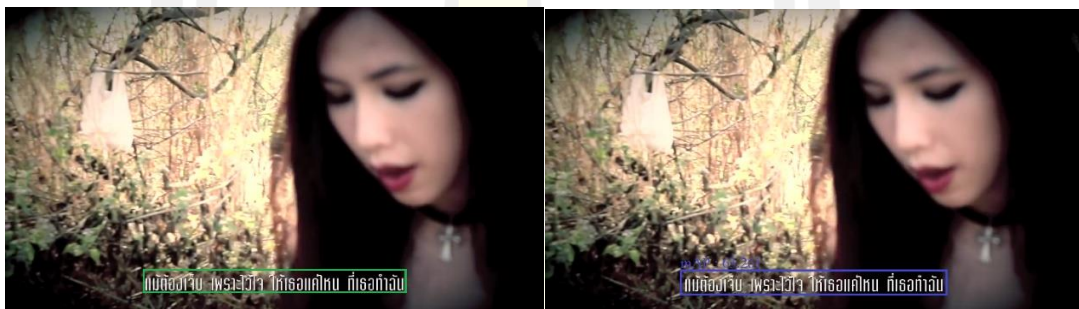
ตารางที่ 4.1 ผลลัพธ์ค่า mAP โดยกำหนดให้ค่า IoU = 0.5

Set	Methods		
	YOLO	RetinaNet	tiny-YOLO
set 1	90.73	90.39	87.00
set 2	85.67	89.56	86.85
set 3	88.84	91.20	88.33
set 4	91.78	92.67	93.33
set 5	87.25	93.32	91.39
mean	88.85	91.43	89.38
SD	2.49	1.56	2.86

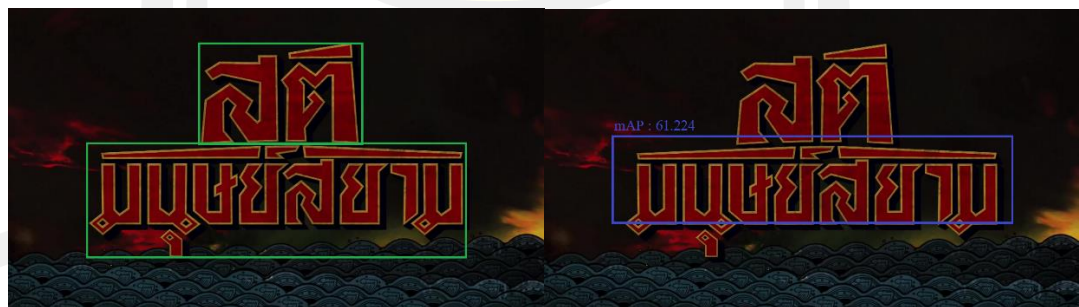
ผลลัพธ์ซึ่งได้จากการทดลองแต่ละวิธีกับชุดของข้อมูลที่แตกต่างกันซึ่งจากการทดลองแสดงให้เห็นว่าผลลัพธ์ที่ดีที่สุดคือ Tiny-YOLO ซึ่งมีค่า mAP 93.33% แต่หากวัดจากค่าเฉลี่ยแล้ววิธีที่ดีที่สุดคือ RetinaNet ซึ่งมีค่า mAP 91.43% ในการตรวจจับคำบรรยายภายใน dataset นี้ตามตารางที่ 4.1(ภาพฝั่งซ้าย) และตัวอย่างผลลัพธ์ในการตรวจจับที่ภาพประกอบที่ 4.1(ภาพฝั่งขวา)



ก)



ข)



ค)

ภาพประกอบที่ 4.1 ภาพผลลัพธ์การตรวจจับตำแหน่งคำบรรยาย รวมไปถึงรอยละที่ถูกต้องและตำแหน่งของกรอบที่ถูกต้องตามภาพ ก) ภาพผลลัพธ์ที่มีกรอบมากกว่า ground truth, ข) ภาพผลลัพธ์ที่มีกรอบเท่ากับ ground truth และ ค) ภาพผลลัพธ์ที่มีกรอบน้อยกว่า ground truth

ผลลัพธ์ค่า CER ซึ่งได้จากการทดลองแต่ละต้นแบบกับชุดข้อมูลที่แตกต่างกัน จากผลการทดลองพบว่า Model A ที่ทดลองโดยใช้ Batch size 64 มีค่า CER ที่น้อยที่สุดคือ 29.93% ตามตารางที่ 4.2 และผลลัพธ์ที่ได้จากการนำสถาปัตยกรรมของ Chamchong et al. มาใช้พบว่า Model 1 ที่ทดลองโดยใช้ Batch size 32 มีค่าน้อยที่สุดคือ 43.28% ตามตารางที่ 4.3 ดังนั้นจึงได้นำข้อมูลค่า CER และเวลาที่ใช้ในการเรียนรู้ (Train) ของ Model A และ Model 1 ที่มีค่า CER น้อยที่สุด มาเปรียบเทียบกันตามตารางที่ 4.4 ซึ่งพบว่า Model A มีผลลัพธ์ที่ดีกว่าโดยมีค่า CER 29.94% และ Model 1 มีค่า CER มากถึง 43.28% แต่ใช้เวลาเพียง 27.09 นาที ซึ่งน้อยกว่าเวลาที่ใช้ในการเรียนรู้ของ Model A 3 เท่า โดยสามารถดูผลลัพธ์ที่ได้จากการรู้จำคำบรรยายด้วย Model A และ Model 1 ได้ที่ภาพประกอบที่ 4.5

ตารางที่ 4.2 ค่า CER ที่ได้จากสถาปัตยกรรมที่ใช้วิธี CNN และ LSTM

Batch Size	CER Value (%)			
	Model A	Model B	Model C	Model D
32	22.61	18.57	17.84	13.50
64	12.37	15.34	19.39	17.12
128	18.38	10.11	9.36	16.53



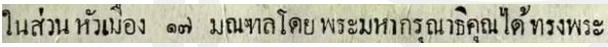


ตารางที่ 4.3 ค่า CER ที่ได้จากการใช้สถาปัตยกรรมตามบทความของ Chamchong et al.

Batch Size	CER Value (%)	
	Model 1	Model 2
32	35.17	14.24
64	68.20	87.22
128	39.22	65.26

ตารางที่ 4.4 ตารางเปรียบเทียบค่า CER และ เวลาในการเรียนรู้ (Train)

Model	Batch Size	CER Value	Training Time (Min.)
C	128	9.36±0.91	115.43
1	32	19.13±1.37	23.12

ตารางที่ 4.5 ผลลัพธ์การทดลองการรู้จำคำบรรยาย

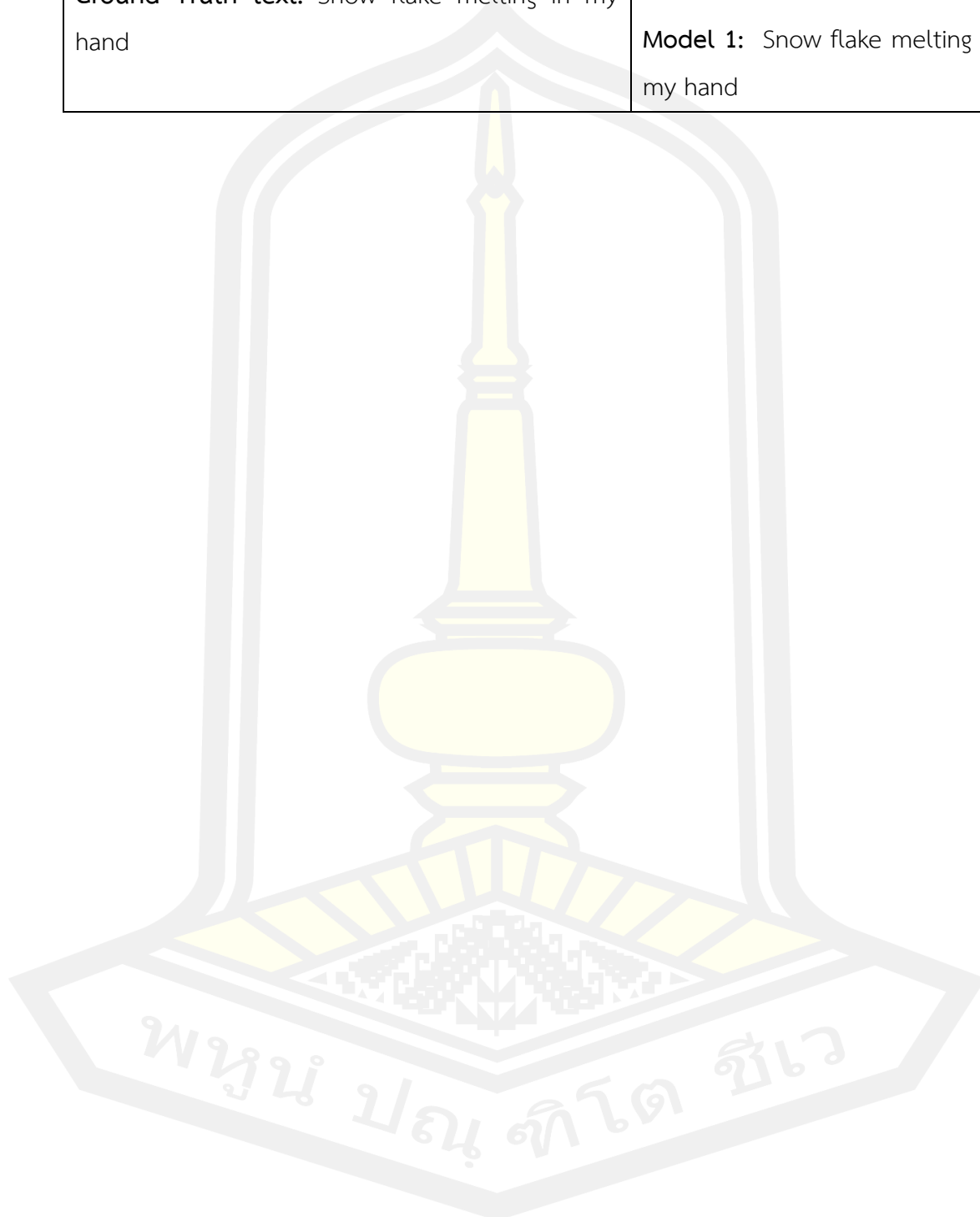
รูปภาพและ Ground truth	ผลลัพธ์การรู้จำ
 <p>Ground Truth text: ให้อาสร้างเราขึ้นใหม่</p>	<p>Model C: ให้อาสร้างเราขึ้นใหม่</p> <p>Model 1: ให้อาสร้างเราขึ้นใหม่</p>
 <p>Ground Truth text: ข้าคนไทย 8 หมื่น ทัวโลก 50 ล้าน!</p>	<p>Model C: ข้าคนไทย 8 หมื่น ทัวโลก 50 ล้าน!</p> <p>Model 1: ข้าคนไทย 2หมื่น กัวโลก 50 ล้านปี</p>
 <p>Ground Truth text: ในส่วนหัวเมือง ๑๗ มณฑลโดย พระมหากษัตริย์คุณได้ทรงพระ</p>	<p>Model C: ในส่วนหัวเมือง ๑๗ มณฑลโดยพระมหากษัตริย์คุณได้ทรงพระ</p> <p>Model 1: ในจันนี้อง ๑๔ พระเศียร ระาทรณ กงยถักรานะ</p>
 <p>Ground Truth text: ทฤษฎีสมคบคิด MH370</p>	<p>Model C: ทฤษฎีสมคบคิด MH370</p> <p>Model 1: ทฤษฎีสมคบคิด MH370</p>
 <p>Ground Truth text: WITHIN YOUR GAZE CONTAINS</p>	<p>Model C: WITHIN YOUR GAZE CONTAINS</p> <p>Model 1: WITHIN YOUR GAZE CONTAINS</p>



Snow flake melting in my hand

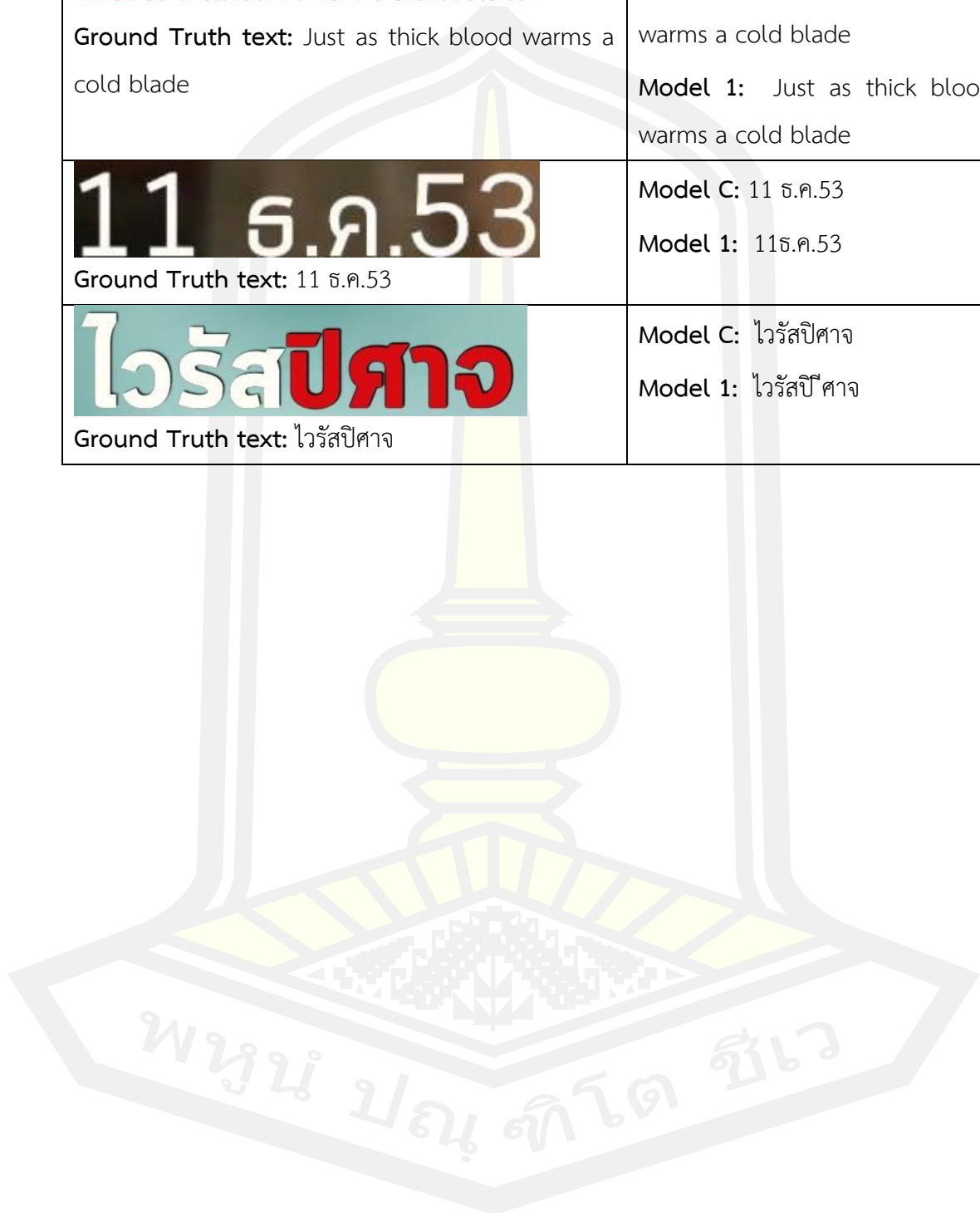
Ground Truth text: Snow flake melting in my hand

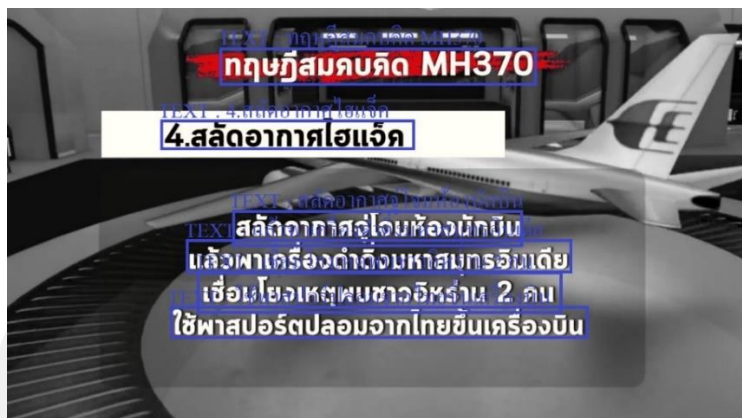
Model C: Snow flake melting in my hand

Model 1: Snow flake melting in my hand



รูปภาพและ Ground truth	ผลลัพธ์การรู้จำ
<p><i>Just as thick blood warms a cold blade</i></p> <p>Ground Truth text: Just as thick blood warms a cold blade</p>	<p>Model C: Just as thick blood warms a cold blade</p> <p>Model 1: Just as thick blood warms a cold blade</p>
 <p>Ground Truth text: 11 ธ.ค.53</p>	<p>Model C: 11 ธ.ค.53</p> <p>Model 1: 11ธ.ค.53</p>
 <p>Ground Truth text: ไวรัสปิตาจ</p>	<p>Model C: ไวรัสปิตาจ</p> <p>Model 1: ไวรัสปีศาจ</p>





ภาพประกอบที่ 4.2 ภาพผลลัพธ์ที่ได้จากการตรวจจับและรู้จำคำบรรยาย

Text: ทฤษฎีสมคบคิด MH370

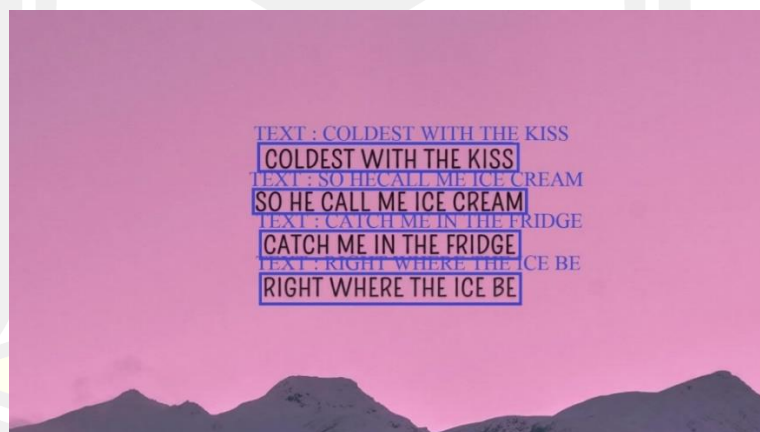
Text: 4.สลัดอากาศไฮแอ็ค

Text: สลัดอากาศดูโอมห้องนักบิน

Text: แล้วพาเครื่องดำดิ่งมหาสมุทรอินเดีย

Text: เชื่อมโยงเหตุพบชาวอิหร่าน 2 คน

Text: ใช้พาสปอร์ตปลอมจากไทยขึ้นเครื่องบิน



ภาพประกอบที่ 4.3 ภาพผลลัพธ์ที่ได้จากการตรวจจับและรู้จำคำบรรยาย ภาษาอังกฤษ

Text: COLDEST WITH THE KISS

Text: SO HE CALL ME ICE CREAM

Text: CATCH ME IN THE FRIDGE

Text: RIGHT WHERE THE ICE BE

บทที่ 5

สรุปและข้อเสนอแนะ

5.1 สรุปผลการวิจัย

งานวิจัยนี้ศึกษาเกี่ยวกับการตรวจจับคำบรรยายโดยใช้ YOLOv3, Tiny-YOLOv3, RetinaNet ซึ่งผลลัพธ์ที่ได้นั้น การเรียนรู้เชิงลึกโดยใช้สถาปัตยกรรมที่เรียกว่า CNN-LSTM ในการรู้จำคำบรรยายจากวิดีโอได้สร้างสถาปัตยกรรม CNN 4 แบบโดยอ้างอิงมาจาก VGG16 และ VGG19 แล้วนำไปเชื่อมกับ Bi-LSTM หรือ Bi-GRU และ ฝึกโดยใช้ CTC Loss

ผลการทดสอบการตรวจจับคำบรรยาย ผลลัพธ์ที่ดีที่สุดโดยวัดจากค่า mAP ที่ใช้ IoU=0.5 พบว่าวิธี YOLOv3 มีค่า mAP 88.85%, Tiny-YOLOv3 มีค่า 89.38% และ RetinaNet มีค่า 91.43% ดังนั้นจึงวิธี RetinaNet จึงเป็นวิธีที่ดีที่สุดในการตรวจจับคำบรรยายของ dataset นี้

ผลการทดสอบการรู้จำคำบรรยาย ผลลัพธ์ที่ดีที่สุดโดยวัดจากค่า CER พบว่า model A มีค่า CER 12.37%, Model B 10.11%, Model C 9.36%, Model D 13.50%, Model 1 17.35%, Model 2 22.87% ซึ่งหากเรียงตามลำดับผลลัพธ์จากดีที่สุดไปแย่สุดได้ดังนี้ Model C, Model B, Model A, Model D, Model 1, และ Model 2 ดังนั้น Model C จึงเป็น Model ที่ดีที่สุดในการรู้จำคำบรรยายของ dataset นี้

ผลลัพธ์การเปรียบเทียบเวลาในการเรียนรู้ของ Model C และ Model 1 ที่มีค่า CER น้อยที่สุดในจำนวนชิ้นเดียวกัน โดยมีจำนวนชั้น CNN 19 และ 4 ชั้นตามลำดับพบว่า Model C ใช้เวลาในการเรียนรู้ 115.43 นาที และ Model 1 ใช้เวลาในการเรียนรู้ 23.12 นาที ซึ่งน้อยกว่า Model C 92.31 นาที หรือ น้อยกว่าประมาณ 3 เท่า

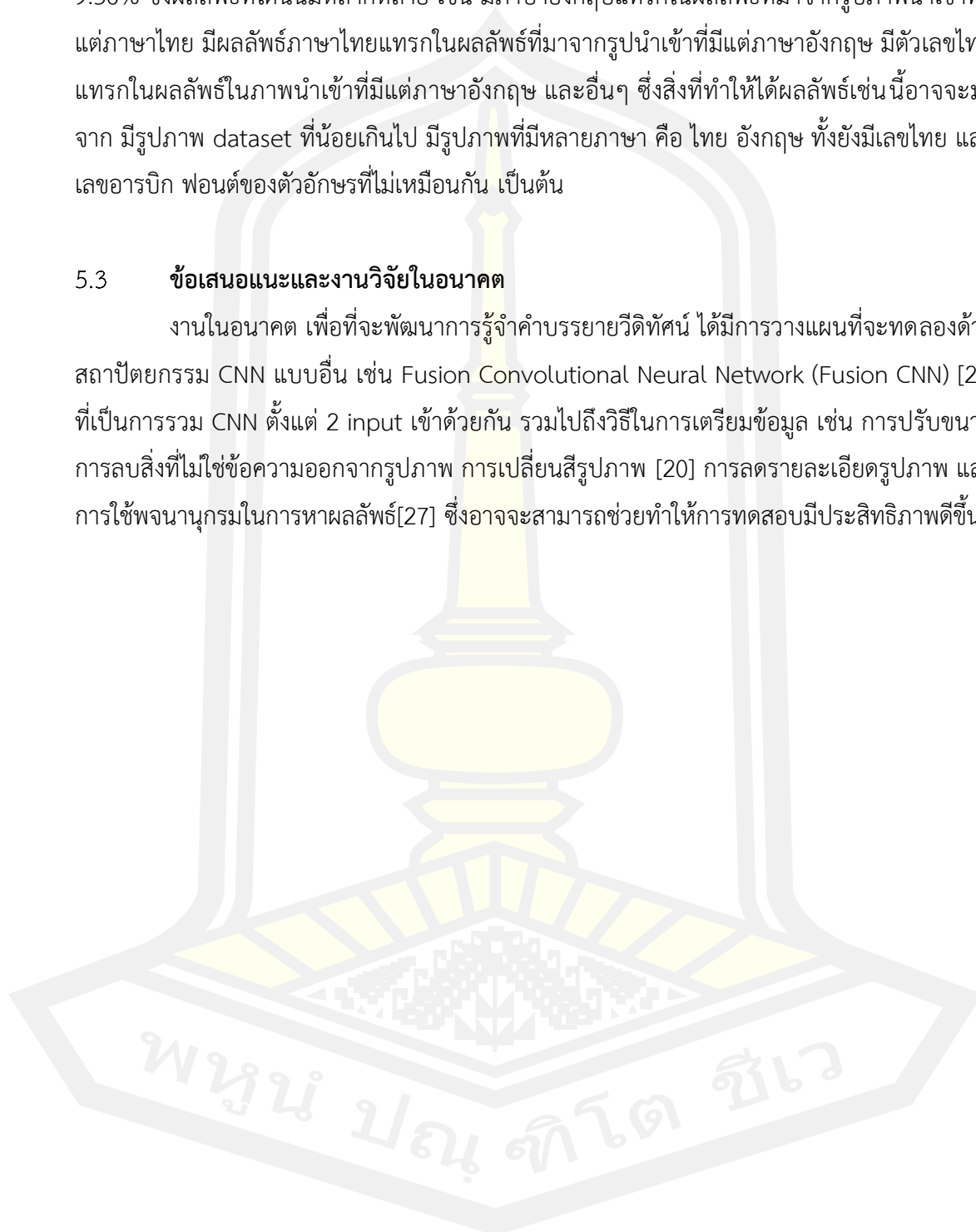
5.2 การอภิปรายผล

จากการทดลองการตรวจจับคำบรรยายวิดีโอพบว่าค่า mAP ที่ดีที่สุดที่ใช้ค่า IoU ที่ 0.5 นั้นคือ 91.43% โดยผลลัพธ์ที่ได้มีส่วนมากจะมีลักษณะเป็นกรอบที่อยู่ตรงกับตำแหน่งของ ground truth มากกว่า 50% แต่ว่ามีบางภาพโดยเฉพาะรูปภาพที่มีคำบรรยายมากกว่า 1 คำบรรยายซึ่งจะทำให้มี ground truth หลายตำแหน่งทำให้ผลลัพธ์ออกมาจะมีจำนวนกรอบมากกว่าคำบรรยายโดยที่กรอบนั้นจะเป็นกรอบใหญ่ซึ่งจะครอบคลุมพื้นที่คำบรรยายทั้งหมดเช่น รูปภาพที่มี 2 คำบรรยาย ทั้ง 2 คำบรรยายนั้นจะมีกรอบที่ครอบพอดีหรือมากกว่า น้อยกว่า เล็กน้อยจากตำแหน่งที่กำหนดใน ground และจะมีกรอบขนาดใหญ่มาครอบ 2 คำบรรยายนั้นเข้าด้วยกันเนื่องจากข้อมูลที่น่าเข้านั้นอาจจะมีพื้นที่ของ ground truth ที่ซ้อนทับกันดังนั้นจึงอาจจะคิดว่าทั้ง 2 คำบรรยายนั้นเป็นคำบรรยายเดียวกัน ทำให้ผลลัพธ์ที่ออกมาจะมีจำนวนกรอบที่มากกว่าจำนวน ground truth

ในส่วนของการทดลองการรู้จำคำบรรยายวิดีโอที่ค้นพบว่าค่าที่ดีที่สุดมีค่า CER 9.36% ซึ่งผลลัพธ์ที่ได้นั้นมีหลากหลาย เช่น มีภาษาอังกฤษแทรกในผลลัพธ์ที่มาจากรูปภาพนำเข้าที่มีแต่ภาษาไทย มีผลลัพธ์ภาษาไทยแทรกในผลลัพธ์ที่มาจากรูปภาพนำเข้าที่มีแต่ภาษาอังกฤษ มีตัวเลขไทยแทรกในผลลัพธ์ในรูปภาพนำเข้าที่มีแต่ภาษาอังกฤษ และอื่นๆ ซึ่งสิ่งที่ทำให้ได้ผลลัพธ์เช่นนี้อาจจะมาจาก มีรูปภาพ dataset ที่น้อยเกินไป มีรูปภาพที่มีหลายภาษา คือ ไทย อังกฤษ ทั้งยังมีเลขไทย และ เลขอารบิก ฟอนต์ของตัวอักษรที่ไม่เหมือนกัน เป็นต้น

5.3 ข้อเสนอแนะและงานวิจัยในอนาคต

งานในอนาคต เพื่อที่จะพัฒนาการรู้จำคำบรรยายวิดีโอ ได้มีการวางแผนที่จะทดลองด้วยสถาปัตยกรรม CNN แบบอื่น เช่น Fusion Convolutional Neural Network (Fusion CNN) [29] ที่เป็นการรวม CNN ตั้งแต่ 2 input เข้าด้วยกัน รวมไปถึงวิธีการเตรียมข้อมูล เช่น การปรับขนาด การลบสิ่งที่ไม่ใช่ข้อความออกจากรูปภาพ การเปลี่ยนสีรูปภาพ [20] การลดรายละเอียดรูปภาพ และการใช้พจนานุกรมในการหาผลลัพธ์[27] ซึ่งอาจจะสามารถช่วยทำให้การทดสอบมีประสิทธิภาพดีขึ้น



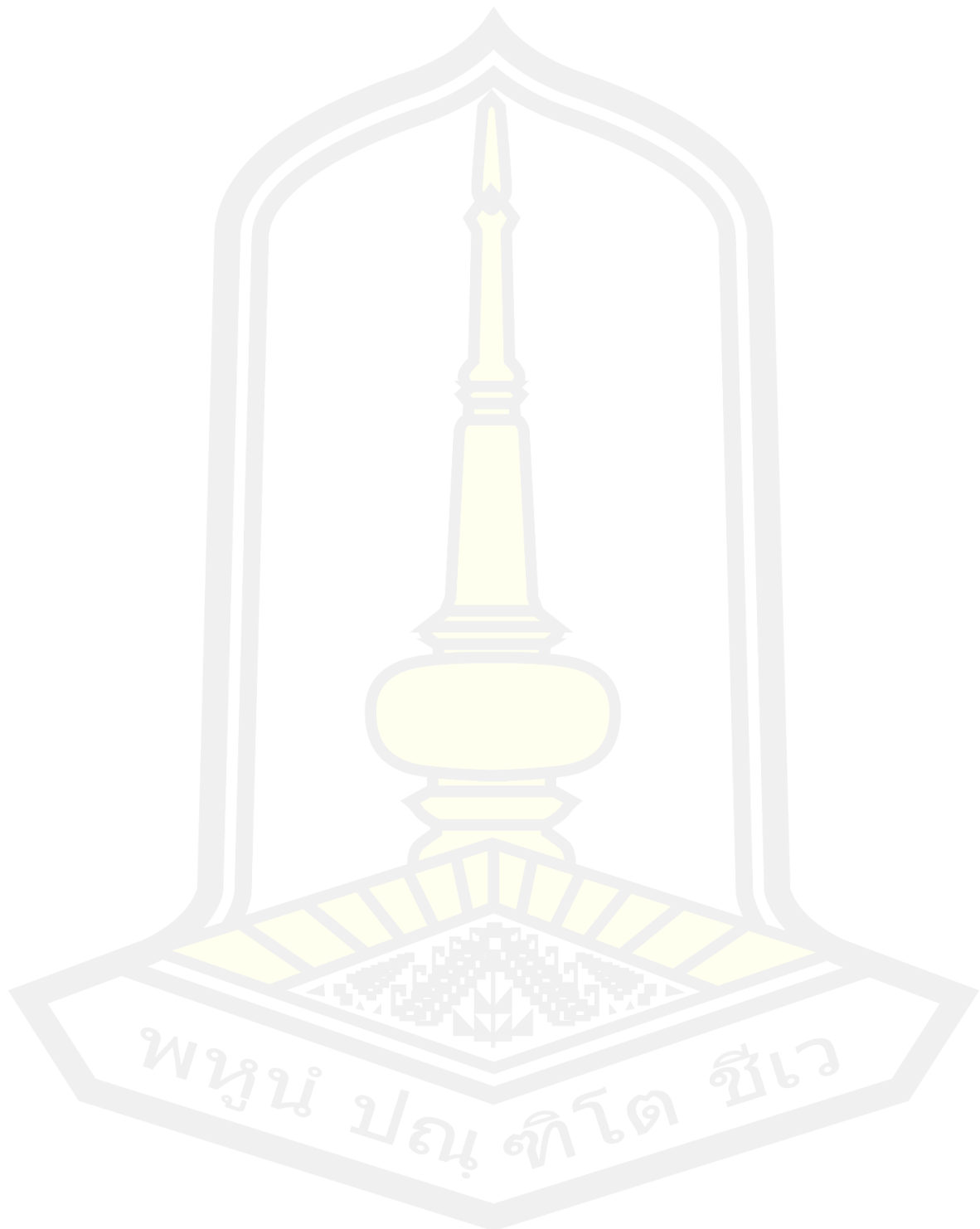
บรรณานุกรม

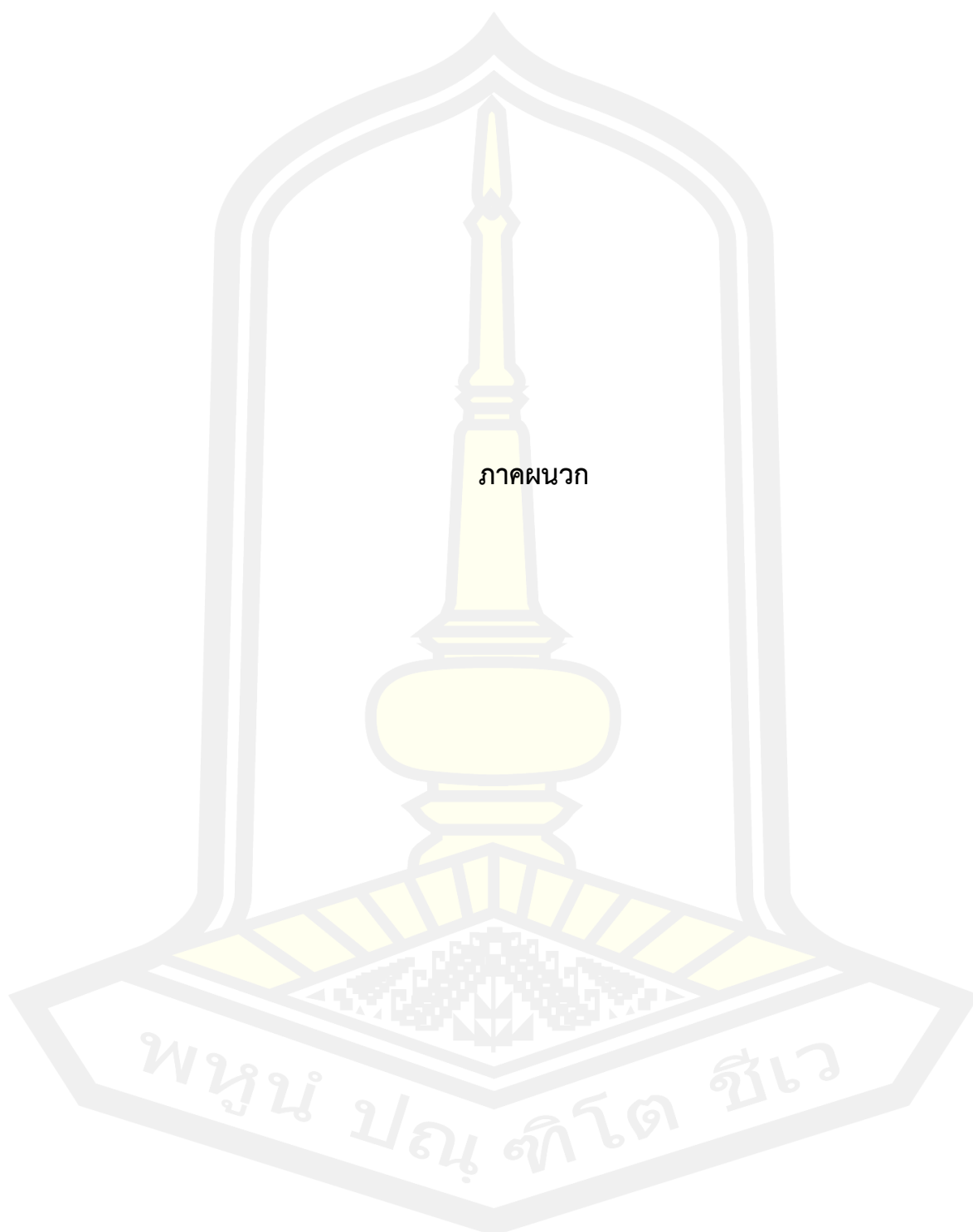
- [1] Z. Zhu, W. Dai, Y. Hu, and J. Li, "Speech Emotion Recognition Model Based on Bi-Gru and Focal Loss," *Pattern Recognition Letters*, vol. 140, pp. 358-365, 2020, doi: <https://doi.org/10.1016/j.patrec.2020.11.009>.
- [2] C. Ma, L. Sun, Z. Zhong, and Q. Huo, "Relatext: Exploiting Visual Relationships for Arbitrary-Shaped Scene Text Detection with Graph Convolutional Networks," *Pattern Recognition*, vol. 111, pp. 107684-107684, 2021, doi: <https://doi.org/10.1016/j.patcog.2020.107684>.
- [3] J. Redmon and A. Farhadi, "Yolov3: An Incremental Improvement," 2018, doi: <https://arxiv.org/abs/1804.02767>.
- [4] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal Loss for Dense Object Detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, pp. 1-1, 2018, doi: 10.1109/tpami.2018.2858826.
- [5] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural computation*, vol. 9, pp. 1735-1780, 1997, doi: 10.1162/neco.1997.9.8.1735.
- [6] K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio, "Learning Phrase Representations Using Rnn Encoder-Decoder for Statistical Machine Translation," 2014, doi: 10.3115/v1/D14-1179.
- [7] S. M. Beitzel, E. C. Jensen, and O. Frieder, "Map," pp. 1691-1692, 2009, doi: 10.1007/978-0-387-39940-9_492.
- [8] P. Wang, R. Sun, H. Zhao, and K. Yu, "A New Word Language Model Evaluation Metric for Character Based Languages," in *Chinese Computational Linguistics and Natural Language Processing Based on Naturally Annotated Big Data*, Berlin, Heidelberg, M. Sun, M. Zhang, D. Lin, and H. Wang, Eds., 2013: Springer Berlin Heidelberg, pp. 315-324, doi: 10.1007/978-3-642-41491-6_29.
- [9] B. P. Woolf, *Chapter 1 - Introduction*. San Francisco: Morgan Kaufmann, 2009, pp. 1-20.
- [10] F. Bre, J. Gimenez, and V. Fachinotti, "Prediction of Wind Pressure Coefficients on Building Surfaces Using Artificial Neural Networks," *Energy and Buildings*, vol.

- 158, 2017, doi: 10.1016/j.enbuild.2017.11.045.
- [11] A. D. Torres, H. Yan, A. H. Aboutaleb, A. Das, L. Duan, and P. Rad, "Patient Facial Emotion Recognition and Sentiment Analysis Using Secure Cloud with Hardware Acceleration," pp. 61-89, 2018, doi: 10.1016/B978-0-12-813314-9.00003-7.
- [12] A. Ram and C. Reyes-Aldasoro, "The Relationship between Fully Connected Layers and Number of Classes for the Analysis of Retinal Images," 2020, doi: <https://arxiv.org/abs/2004.03624>.
- [13] Y. Xu *et al.*, "End-to-End Subtitle Detection and Recognition for Videos in East Asian Languages Via Cnn Ensemble," *Signal Processing: Image Communication*, vol. 60, pp. 131-143, 02/01 2018, doi: <https://doi.org/10.1016/j.image.2017.09.013>.
- [14] He, C.-W. Huang, L. Wei, L. Li, and G. Anfu, "Tf-Yolo: An Improved Incremental Network for Real-Time Object Detection," *Applied Sciences*, vol. 9, pp. 3225-3225, 2019, doi: 10.3390/app9163225.
- [15] T. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature Pyramid Networks for Object Detection," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 21-26 July 2017, pp. 936-944, doi: 10.1109/CVPR.2017.106.
- [16] J. Gan, W. Wang, and K. Lu, "In-Air Handwritten Chinese Text Recognition with Temporal Convolutional Recurrent Network," *Pattern Recognition*, vol. 97, p. 107025, 2020/01/01/ 2020, doi: <https://doi.org/10.1016/j.patcog.2019.107025>.
- [17] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, *Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks*. 2006, pp. 369-376.
- [18] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *arXiv 1409.1556*, 09/04 2014.
- [19] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized Intersection over Union: A Metric and a Loss for Bounding Box Regression," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 15-20 June 2019, pp. 658-666, doi: 10.1109/CVPR.2019.00075.

- [20] W. He, X.-Y. Zhang, F. Yin, Z. Luo, J.-M. Ogier, and C.-L. Liu, "Realtime Multi-Scale Scene Text Detection with Scale-Based Region Proposal Network," *Pattern Recognition*, vol. 98, pp. 107026-107026, 2020, doi: <https://doi.org/10.1016/j.patcog.2019.107026>.
- [21] Y. Wang, L. Peng, and S. Wang, "A Multi-Stage Method for Chinese Text Detection in News Videos," *Procedia Computer Science*, vol. 96, pp. 1409-1417, 2016, doi: <https://doi.org/10.1016/j.procs.2016.08.186>.
- [22] Y. Zhu and J. Du, "Textmountain: Accurate Scene Text Detection Via Instance Segmentation," *Pattern Recognition*, vol. 110, pp. 107336-107336, 2021, doi: <https://doi.org/10.1016/j.patcog.2020.107336>.
- [23] M. Huang, C. Lan, W. Huang, and Y. Tao, "Natural Scene Text Detection Based on Multiscale Connectionist Text Proposal Network," *The Journal of Engineering*, vol. 2020, no. 13, pp. 326-329, 2020, doi: 10.1049/joe.2019.1154.
- [24] H. Yan and X. Xu, "End-to-End Video Subtitle Recognition Via a Deep Residual Neural Network," *Pattern Recognition Letters*, vol. 131, pp. 368-375, 2020, doi: <https://doi.org/10.1016/j.patrec.2020.01.019>.
- [25] S. K. Jemni, Y. Kessentini, and S. Kanoun, "Out of Vocabulary Word Detection and Recovery in Arabic Handwritten Text Recognition," *Pattern Recognition*, vol. 93, pp. 507-520, 2019, doi: <https://doi.org/10.1016/j.patcog.2019.05.003>.
- [26] J. Gan, W. Wang, and K. Lu, "In-Air Handwritten Chinese Text Recognition with Temporal Convolutional Recurrent Network," *Pattern Recognition*, vol. 97, pp. 107025-107025, 2020, doi: <https://doi.org/10.1016/j.patcog.2019.107025>.
- [27] J. Zhang, C. Luo, L. Jin, T. Wang, Z. Li, and W. Zhou, "Sahan: Scale-Aware Hierarchical Attention Network for Scene Text Recognition," *Pattern Recognition Letters*, vol. 136, pp. 205-211, 2020, doi: <https://doi.org/10.1016/j.patrec.2020.06.009>.
- [28] R. Chamchong, W. Gao, and M. D. McDonnell, "Thai Handwritten Recognition on Text Block-Based from Thai Archive Manuscripts," 2019, pp. 1346-1351, doi: 10.1109/ICDAR.2019.00217.
- [29] Y. Lavinia, H. H. Vo, and A. Verma, "Fusion Based Deep Cnn for Improved Large-Scale Image Action Recognition," in *IEEE International Symposium on*

Multimedia (ISM), 11-13 Dec. 2016, pp. 609-614, doi: 10.1109/ISM.2016.0131.





ภาคผนวก

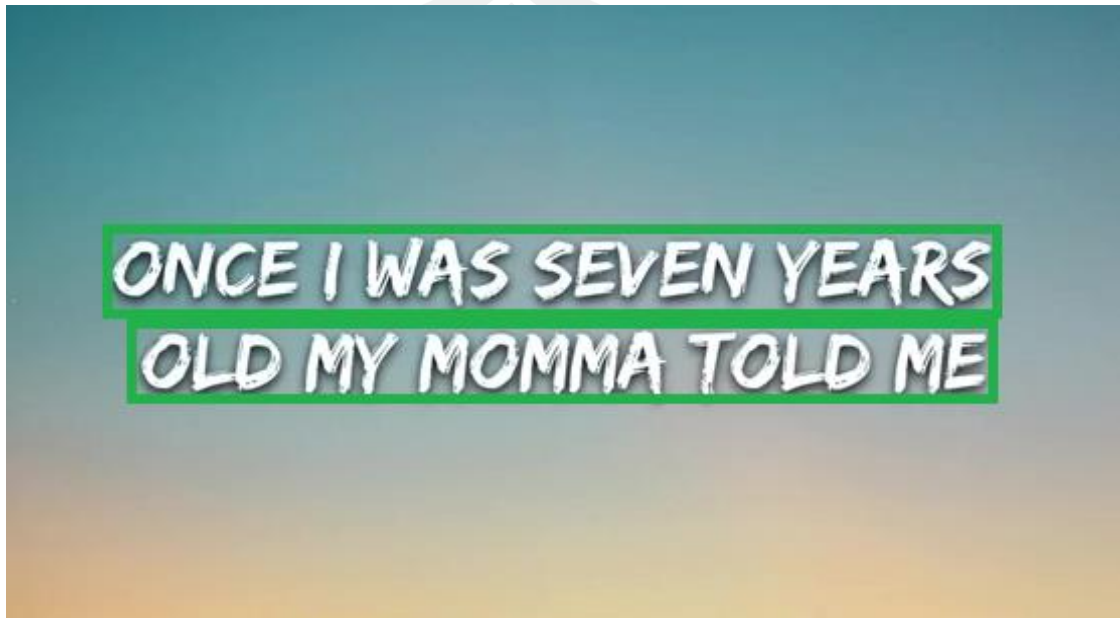
พหุบัณฑิต วิไล



ภาคผนวก ก

ตัวอย่างรูปภาพที่ตัดจากวีดิทัศน์และคำตอบ (Ground truth) แสดงบริเวณที่เป็นคำบรรยาย

สามารถดาวน์โหลดได้จากลิงก์: <https://drive.google.com/drive/folders/1w212u-exPzbS-IZTM5s8VDXmX69GFluj?usp=sharing>





su shuo ji fan qing hua kuo

of many love words that I want to say
let the words come out a sound



LIKE A BULLET
IN MY CHEST
YOU'RE WRITTEN
BOUND AND ETCHED

EQUILANORA





네 곁에서 멀어지지 않아

ne gyeoteseo meoreojji anha

I'm not getting far away by your side



SHÒU HUÒ MÍ DÌ DE DÀI JIÀ

是收获谜底的代价

[KORIAN: SHONHUMIKODAEJINGEORANWADAEWISHOHWATAM]

พหุจน์ ปณฺ ทิโต ชีเว



หนังสือพิมพ์สมัยนั้น ได้ลงเรื่องราวของ

พหุจน์ ปณฺ ทิโต ชีเว



พหุบัน ปณุ ทิโต ชีเว

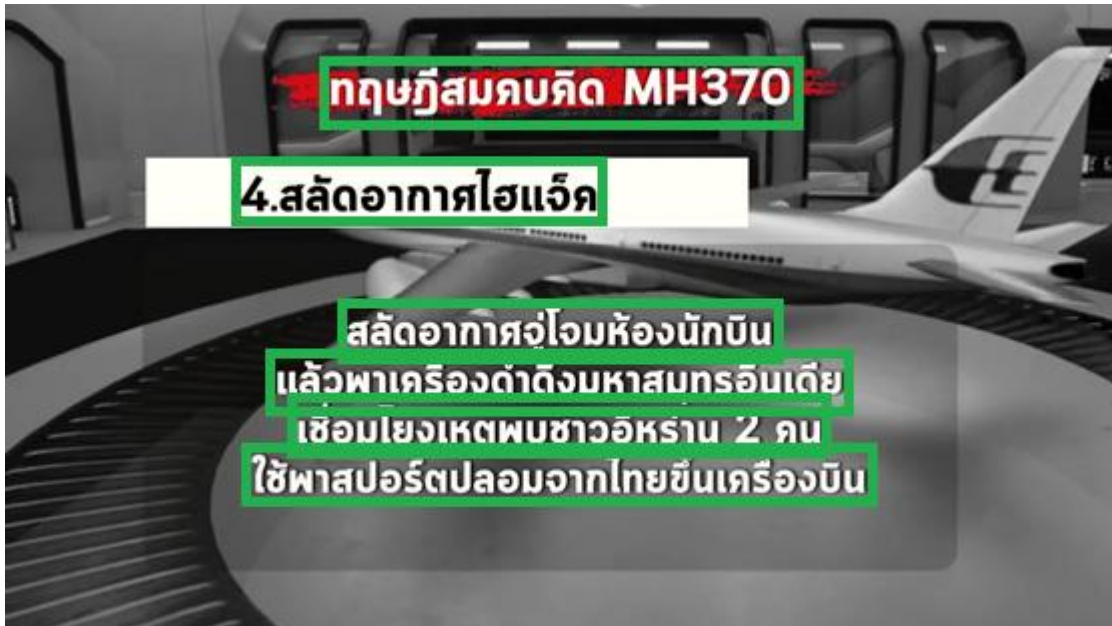


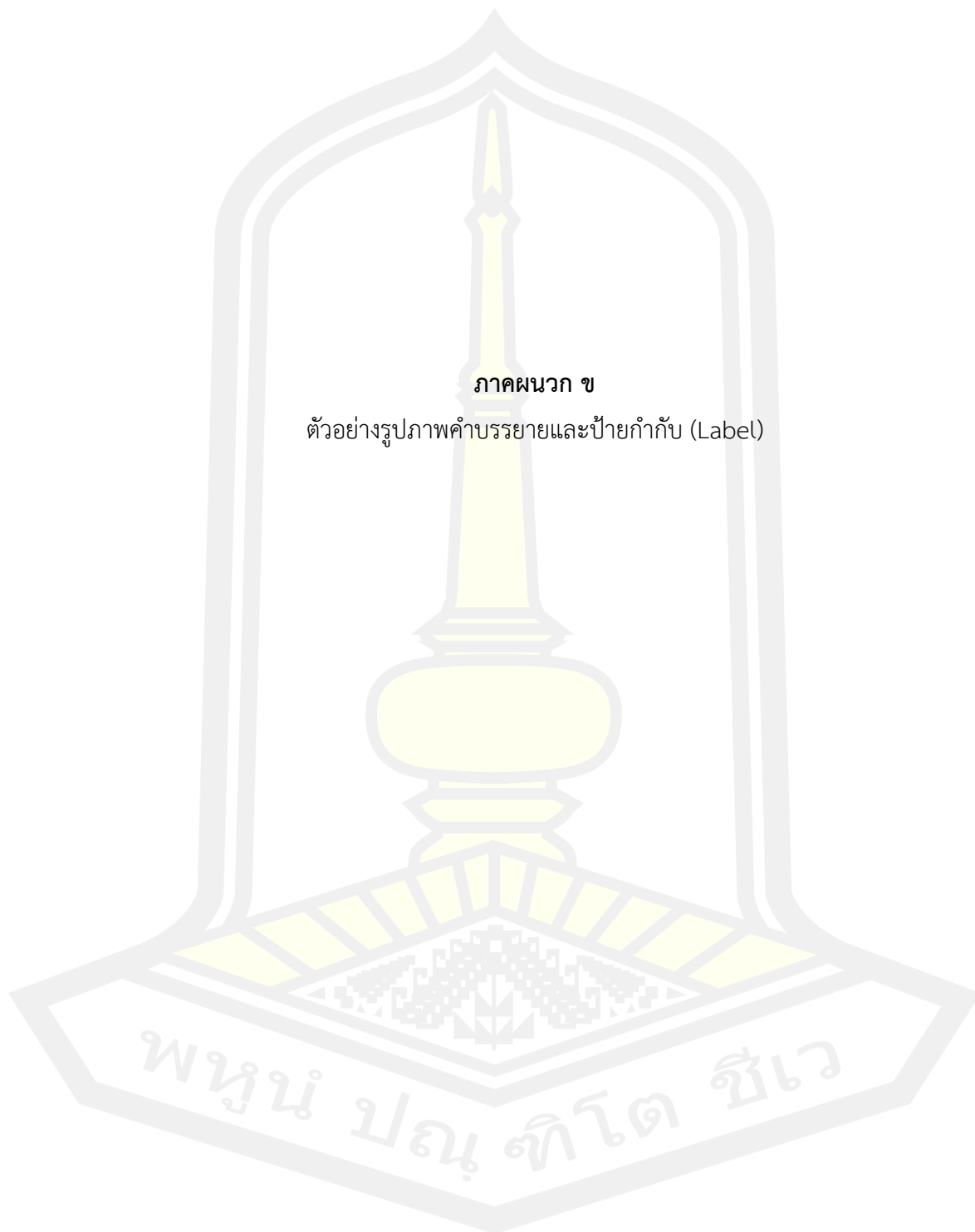


 NocNoc.com

เข้ามาแต่งบ้านให้สุดอย่างฟินที่

พหุบัน ปณู ทิโต ชีเว





ภาคผนวก ข

ตัวอย่างรูปภาพคำบรรยายและป้ายกำกับ (Label)

สามารถดาวน์โหลดได้จากลิงก์: <https://drive.google.com/drive/folders/1w212u-exPzbS-IZTM5s8VDXMx69GFluj?usp=sharing>

ตัวอย่างภาษาไทย

ผู้รักชากฎหมาย

Label: ผู้รักชากฎหมาย

ต้องลงมือเด็ดขาด

Label: ต้องลงมือเด็ดขาด

ทอนซิลอักเสบ

Label: ทอนซิลอักเสบ

ไข່สูง ปวดศีรษะ

Label: ไข່สูง ปวดศีรษะ

จำเลือดตามผิวหนัง

Label: จำเลือดตามผิวหนัง

ในส่วน หัวเมือง ๑๗ มณฑลโดยพระมหากษัตริย์คุณได้ทรงพระ

Label: ในส่วน หัวเมือง ๑๗ มณฑลโดยพระมหากษัตริย์คุณได้ทรงพระ

พระราชอาณาจักร ฟังสงบลงเมื่อเดือนมีนาคมที่ล่วงมานั้น

Label: พระราชอาณาจักร ฟังสงบลงเมื่อเดือนมีนาคมที่ล่วงมานั้น

กรุณาโปรดเกล้าฯ ให้กรมสาธารณสุขจัดส่งแพทย์ ส่งยาและคำ

Label: กรุณาโปรดเกล้าฯ ให้กรมสาธารณสุขจัดส่งแพทย์ ส่งยาและคำ

ถ้าเป็นภาษาหนึ่ง **“ผมโดนเป็นสิบลๆ เทคเลย”**

Label: ถ้าเป็นภาษาหนึ่ง “ผมโดนเป็นสิบลๆ เทคเลย”

ถ้าอย่างใดอย่างหนึ่งผิดก็ต้องเริ่มใหม่

Label: ถ้าอย่างใดอย่างหนึ่งผิดก็ต้องเริ่มใหม่

ตอนนั้นผมเรียนเพาะช่างไปเจอสาวคนหนึ่งบนรถเมล์

Label: ตอนนั้นผมเรียนเพาะช่างไปเจอสาวคนหนึ่งบนรถเมล์

ต่อมา ศาลชั้นต้นพิพากษาประหารชีวิต****

Label: ต่อมา ศาลชั้นต้นพิพากษาประหารชีวิต

อาจตกอยู่ในสภาพโคม่่าไม่ซ่าไม่นาน****

Label: อาจตกอยู่ในสภาพโคม่่าไม่ซ่าไม่นาน

พหุบุ ปณ จิต โต ชีเว

ตัวอย่างภาษาอังกฤษ

COMMANDER

Label: COMMANDER

I REALLY LOVE

Label: I REALLY LOVE

CAUSE ONLY THOSE

Label: CAUSE ONLY THOSE

WILL EVER REALLY KNOW ME

Label: WILL EVER REALLY KNOW ME

Alone, waiting a thousand years just to return

Label: Alone, waiting a thousand years just to return

Dreaming of sentiments, feelings spreading a long distance

Label: Dreaming of sentiments, feelings spreading a long distance

Deep sadness, tears flow because of this love story

Label: Deep sadness, tears flow because of this love story

Ju Jing Yi

Label: Ju Jing Yi

Ancient Painting

Label: Ancient Painting

Having been through spring and fall, winter and summer

Label: Having been through spring and fall, winter and summer

HOW COULD I JUST GLEEFULLY COMPLY PLACIDLY?

Label: HOW COULD I JUST GLEEFULLY COMPLY PLACIDLY?

YOU CAN ONLY LOSE WHAT YOU CLING TO

Label: YOU CAN ONLY LOSE WHAT YOU CLING TO

INSTANTLY WHEN I THINK OF YOU — I JUST

Label: INSTANTLY WHEN I THINK OF YOU – I JUST

GET IT FREE LIKE WILLY

Label: GET IT FREE LIKE WILLY

SNOW CONE CHILLY

Label: SNOW CONE CHILLY

IN THE JEANS LIKE BILLIE

Label: IN THE JEANS LIKE BILLIE

Throughout the entire vacation, reading One Hundred Years of Solitude

Label: Throughout the entire vacation, reading One Hundred Years of Solitude

Letting the tears stream, diluting her tipsiness

Label: Letting the tears stream, diluting her tipsiness

Playing the songs she loved as a girl is enough to make her dance

Label: Playing the songs she loved as a girl is enough to make her dance

TOAST TO THE ONES HERE TODAY

Label: TOAST TO THE ONES HERE TODAY

OF EVERYTHING WE'VE BEEN THROUGH

Label: OF EVERYTHING WE'VE BEEN THROUGH

TOAST TO THE ONES THAT WE LOST ON THE WAY

Label: TOAST TO THE ONES THAT WE LOST ON THE WAY

I DON'T GIVE A FUCK ABOUT YOU ANYWAYS

Label: I DON'T GIVE A FUCK ABOUT YOU ANYWAYS



ตัวอย่างที่มีการผสมภาษาไทย ภาษาอังกฤษอังกฤษ และ ตัวเลข

11 ธ.ค.53

Label: 11 ธ.ค.53

คดีนี้มีพยานเสียชีวิตกะทันหันรวม 3 คน

Label: คดีนี้มีพยานเสียชีวิตกะทันหันรวม 3 คน

ปี 2016 โลกได้พบชิ้นส่วนปีกของเครื่อง MH370

Label: ปี 2016 โลกได้พบชิ้นส่วนปีกของเครื่อง MH370

MH370 ไม่มีคนคุมระหว่างตก เพียงแต่บินไปเรื่อยๆ

Label: MH370



ประวัติผู้เขียน

ชื่อ	ชนดล สิงขรอาสน์
วันเกิด	13 เมษายน 2541
สถานที่เกิด	จังหวัดมหาสารคาม
สถานที่อยู่ปัจจุบัน	5 ถนนริมคลองสมถวิล ตำบลตลาด อำเภอเมือง จังหวัดมหาสารคาม 44000
ตำแหน่งหน้าที่การงาน	นักวิชาการคอมพิวเตอร์
สถานที่ทำงานปัจจุบัน	กองแผนงาน มหาวิทยาลัยมหาสารคาม
ประวัติการศึกษา	พ.ศ. 2558 โรงเรียนสาธิตมหาวิทยาลัยมหาสารคาม (ฝ่ายมัธยม) พ.ศ. 2562 วิทยาศาสตรบัณฑิต (วท.บ.) สาขาเทคโนโลยีสารสนเทศ มหาวิทยาลัยมหาสารคาม พ.ศ. 2564 วิทยาศาสตรมหาบัณฑิต (วท.ม.) สาขาเทคโนโลยีสารสนเทศ มหาวิทยาลัยมหาสารคาม

พูนุ ปณุกิตโต ชีวะ