

บทที่ 3

วิธีดำเนินการวิจัย

สำหรับวิธีการดำเนินการวิจัยการพัฒนาระบบการรู้จำท่าทางภาษามือไทยด้วยโครงข่ายประสาทเทียมแบบวนกลับ นั้นสามารถแบ่งออกเป็น 5 ส่วนดังนี้

- 3.1 การเตรียมข้อมูล
- 3.2 การฝึกฝนโมเดล
- 3.3 การวัดประสิทธิภาพโมเดล
- 3.4 การเปรียบเทียบประสิทธิภาพโมเดล
- 3.5 การทดสอบโมเดล

3.1 การเตรียมข้อมูล

3.1.1 การรวบรวมข้อมูล

ในการรวบรวมข้อมูล สำหรับการสร้าง TSL10 (dataset ภาษามือไทย 10 ท่า) ผู้วิจัยต้องการวิดีโอท่าภาษามือที่ใช้ในชีวิตประจำวันของผู้พิการทางการได้ยินและการสื่อความหมาย เป็นจำนวน 10 คำ ซึ่งเป็นท่าที่นำมาจาก เว็บไซต์ highlight.kapook.com ที่เนื้อหาเกี่ยวกับการแนะนำภาษามือเบื้องต้นสำหรับใช้ในชีวิตประจำวัน ผู้วิจัยได้มีการออกหนังสือขอความอนุเคราะห์จากศูนย์บริการสนับสนุนการนักศึกษาพิการระดับอุดมศึกษา (DSS) ประจำมหาวิทยาลัยราชภัฏสกลนครเพื่อเก็บข้อมูลสำหรับการเทรนโมเดลสำหรับการรู้จำภาษามือไทยจากทั้งผู้เชี่ยวชาญภาษามือและผู้พิการที่ใช้ภาษามือเป็นหลักในการสื่อสาร โดยผู้วิจัยจะทำเป็นวิดีโอ 85 วิดีโอต่อ 1 คำ และใน 1 วิดีโออัตราเฟรมต่อวินาทีที่ 30 FPS ขนาดของวิดีโอคือ 640 x 480

ตารางที่ 3.1 คำศัพท์ภาษามือที่ใช้ในโครงการ

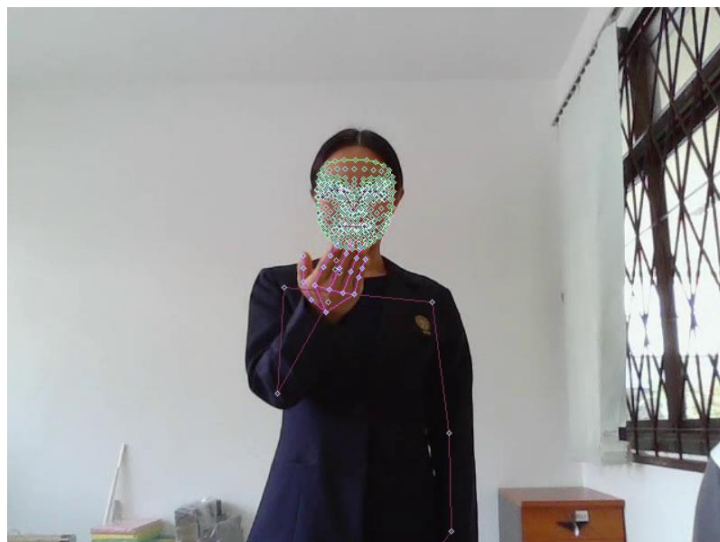
คำภาษาไทย	คำภาษาอังกฤษ	ความหมาย
ขอบคุณ	Thank You	กล่าวแสดงความรู้สึกถึงบุญคุณหรือกล่าวเมื่อได้รับความช่วยเหลือ
ขอโทษ	Sorry	ขอภัยเมื่อได้ทำผิดพลาดอย่างใดอย่างหนึ่ง
ไม่เป็นไร	That is OK	คำแสดงความรู้สึกที่ไม่ได้ถือโทษหรือโกรธเคืองใด ๆ เพื่อให้ผู้ฟังรู้สึกดีขึ้นหรือไม่ต้องรู้สึกผิด
สบายดี	Fine	สภาวะปกติของทั้งร่างกายและจิตใจ ร่างกายไม่เจ็บป่วย รวมทั้งอารมณ์ดี มีความสุข ไม่มีอะไรให้กังวล
ชอบ	Like	พอใจ แสดงอาการพึงพอใจ
รัก	Love	มีใจผูกพันอย่างมาก
ไม่สบาย	Sick	สภาวะที่ร่างกายและจิตใจไม่ปกติ หรือเกิดอาการป่วย
สวัสดี	Hello	ใช้สำหรับการทักทายผู้คน
ฉัน	IAm	ใช้สำหรับการเรียกแทนตัวเอง
คุณ	You	ใช้สำหรับเรียกแทนผู้ที่เราพูดด้วย



ภาพที่ 3.1 ตัวอย่างภาษามือไทย ‘สวัสดี’ จากผู้เชี่ยวชาญภาษามือไทย

3.1.2 การสกัดลักษณะเด่นของข้อมูล

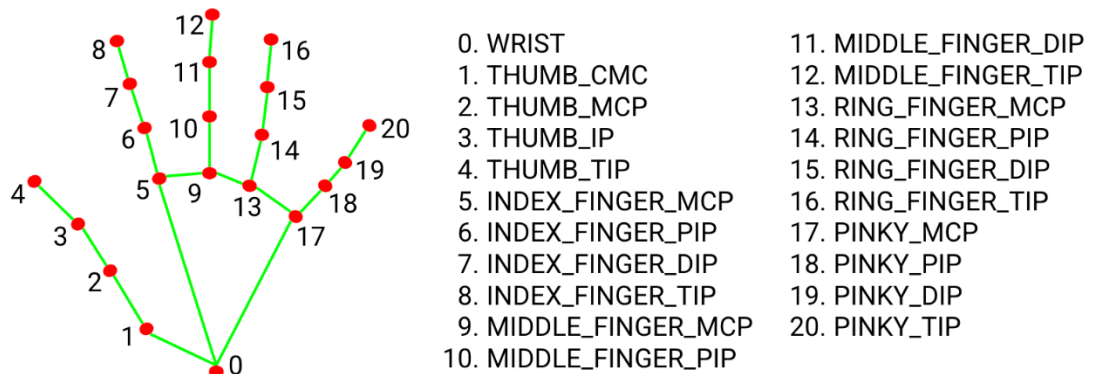
ภาษามือนั้นขึ้นอยู่กับการใช้มือและท่าทาง การนำวิดีโอที่เป็นภาษามือมาใช้ในการเทรนโมเดลนั้นจึงเป็นเรื่องยาก ผู้วิจัยจึงได้ใช้เครื่องมือ MediaPipe ที่เป็น Framework มาใช้ในการแก้ปัญหา ซึ่งวิธีการคือการใช้ MediaPipe ในการ Keypoints ขึ้นตามจุดต่าง ๆ ของร่างกายเป็นค่ามิติ X, Y, Z ของหน้า, มือและท่าทางรูปภาพที่ 3.2



ภาพที่ 3.2 การใช้ MediaPipe ในการ Keypoints

ในมือแต่ละข้างนั้น MediaPipe จะสกัดออกมาได้ 21 Keypoints ซึ่ง Keypoint จะถูกคำนวณแบบ 3 มิติ X, Y, Z ของมือทั้งสองข้าง โดยจะได้ Keypoints จากการสกัดจากมือดังนี้

Keypoints in hand x Three dimensions x No. of hands = $(21 \times 3 \times 2) = 126$ Keypoints ดังภาพที่ 3.3

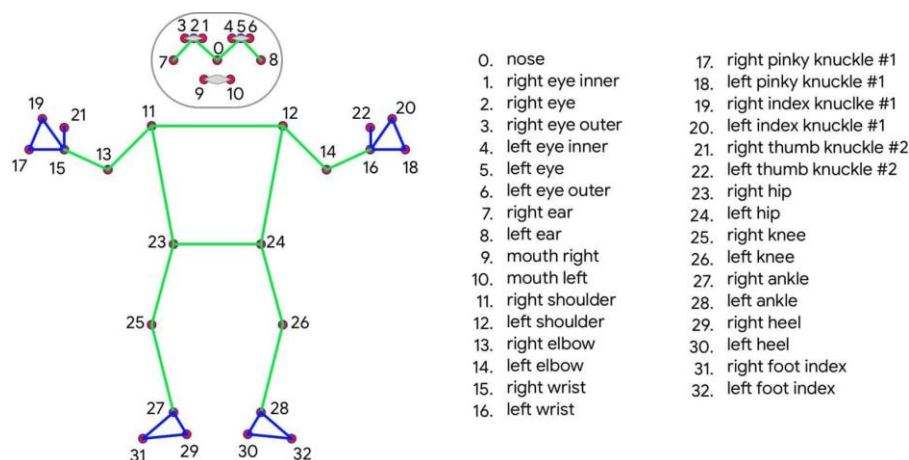


ภาพที่ 3.3 ลำดับและป้ายกำกับ Keypoints ของมือใน MediaPipe

ที่มา : MediaPipe (2023: Online)

ในส่วนของการท่าทางนั้น MediaPipe จะสกัดออกมาได้ 33 Keypoints คำนวณแบบ 3 มิติ X, Y, Z และเพิ่มค่า Visibility เข้าไปซึ่งเป็นค่าที่จะระบุว่าจุดนั้นมองเห็นหรือซ่อนอยู่ (ที่ถูกปิดโดยจุดอื่นของร่างกาย) บนเฟรมดังนั้นจะได้ค่า Keypoints ดังนี้

Keypoints in pose x (Three dimensions + Visibility) = $(33 + (33 + 1)) = 132$ Keypoints ดังภาพที่ 3.4



ภาพที่ 3.4 ลำดับและป้ายกำกับ Keypoints ของท่าทางใน MediaPipe

ที่มา : MediaPipe (2023: Online)

สำหรับหน้านั้น Mediapipe สกัดออกมาได้ 468 Keypoints ได้แก่ รูปทรงรอบหน้าและหน้า, ตา, ปากและคิ้ว ซึ่งคำนวณค่า 3 มิติ X, Y, Z ได้ดังนี้

Keypoints in face x Three dimensions = $(468 \times 3) = 1404$ Keypoints ดังภาพที่ 3.5



ภาพที่ 3.5 Keypoints บนหน้า

ดังนั้นเมื่อรวม Keypoint ทั้งหมดเข้าด้วยกันไม่ว่าจะเป็นจาก หน้า ท่าทางและมือจะสามารถคำนวณได้ดังนี้

Keypoints in hands + in pose + inface = $(126 + 132 + 1404) = 1662$ Keypoints

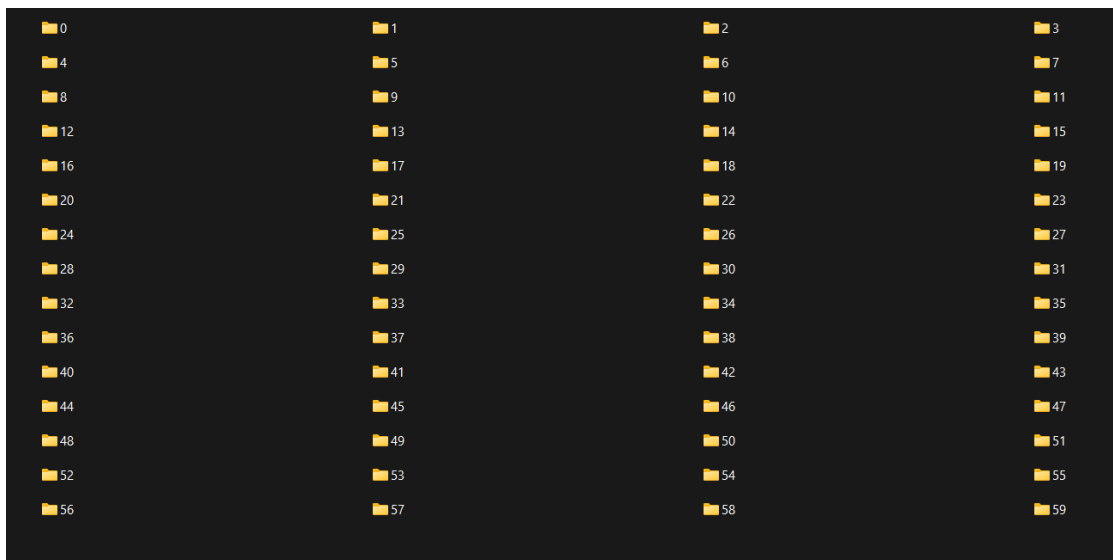
3.1.3 การเตรียมไฟล์

เมื่อสามารถสร้าง Keypoints เสร็จขั้นตอนต่อไปคือการนำผลของค่า Keypoints ของแต่ละจุดของร่างกายเขียนเป็น .npy ไฟล์ ซึ่งมีขั้นตอนดังนี้

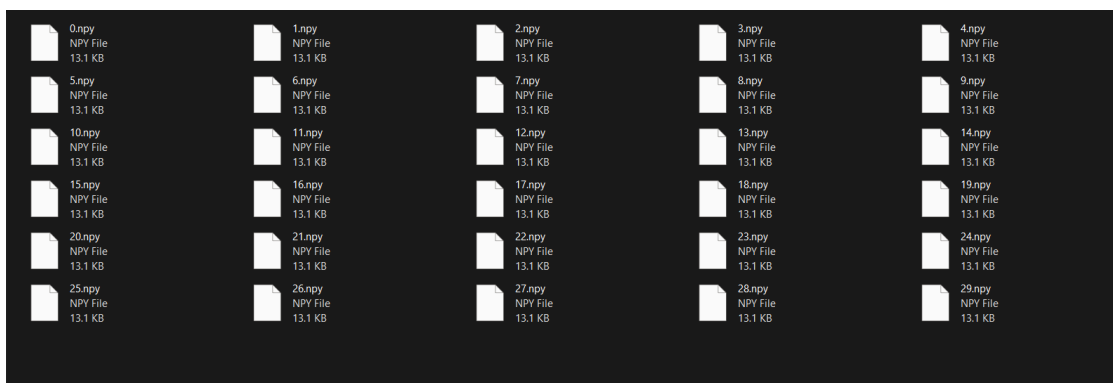
1. สร้างโฟลเดอร์สำหรับเก็บ Datasets
2. ในโฟลเดอร์ Datasets มี โฟลเดอร์ที่เป็นชื่อท่าภาษามือ ดังภาพที่ 3.6
3. ในโฟลเดอร์ที่เป็นชื่อท่าภาษามือจะมีโฟลเดอร์สำหรับเก็บวิดีโอท่าภาษามือ 85 วิดีโอ โดยแยกเป็นโฟลเดอร์ละ 1 วิดีโอ ดังภาพที่ 3.7
4. ในโฟลเดอร์เก็บวิดีโอท่าภาษามือจะมีไฟล์ .npy 30 ไฟล์ ซึ่ง 1 ไฟล์ จะเก็บค่าที่ได้จากการสกัด Keypoints จาก Mediapipe X, Y, Z ใน 1 เฟรม ดังภาพที่ 3.8



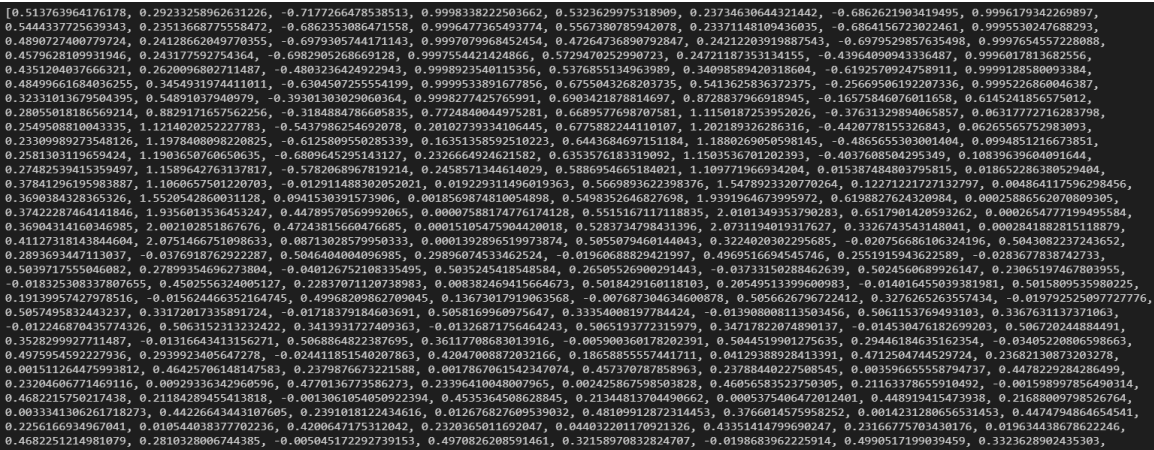
ภาพที่ 3.6 โฟลเดอร์ชื่อท่าภาษามือ



ภาพที่ 3.7 โฟลเดอร์ 60 โฟลเดอร์สำหรับเก็บ .npy ไฟล์



ภาพที่ 3.8 ไฟล์ .npy 30 ไฟล์ ใน 1 โฟลเดอร์วิดีโอ



ภาพที่ 3.9 ไฟล์ .npy ที่เก็บค่า X, Y, Z ของ Keypoints

3.2 การฝึกฝนโมเดล

ผู้วิจัยได้ใช้โมเดลในการเทรนทั้งหมด 3 โมเดลได้แก่ LSTM, GRU, BiLSTM ในงานวิจัยครั้งนี้ ซึ่งเป็นโมเดลของ Recurrent Neurons Networks (RNN)

Number of Nodes คือ จำนวนของ Input Node ซึ่งผู้วิจัยกำหนดขั้นต่ำไว้ 64 จนถึง 256 Activation คือตัวฟังก์ชันที่ใช้ในการรับผลรวมจากการประมวลผลทั้งหมดจากทุก Input Node เข้ามาพิจารณาตามกลไกการคำนวณของ Activation Function นั้น ๆ แล้วส่งต่อไปเป็น Output ซึ่งในงานวิจัยนี้ได้เลือกใช้ 2 ตัว คือ Rectified Linear Unit (ReLU) และ Softmax

Optimizer คือ อัลกอริทึมการเพิ่มประสิทธิภาพ (Optimizer) ทำหน้าที่เป็นกลไกการปรับปรุงค่าน้ำหนักของตัวแปรต้นต่าง ๆ รวมถึงค่าความคลาดเคลื่อน (Bias) ในงานวิจัยนี้ได้เลือกใช้ Optimizer ได้แก่ Adagrad, Adamax, Adam or RMSprop ดังตารางที่ 3.2.1

ตารางที่ 3.2 พารามิเตอร์ของเลเยอร์โมเดล

Parameters	Value
RNN Model	GRU, LSTM, BiLSTM
Number of Nodes	Between (64, 256)
Activation	‘Relu’ or ‘Softmax’
Optimizer	‘Adagrad’, ‘Adamax’, ‘Adam’ or ‘RMSprop’

3.3 การวัดประสิทธิภาพโมเดล

การวัดประสิทธิภาพของโมเดล ผู้วิจัยได้ใช้ตัวชี้วัดคือค่า Accuracy หรือก็คือค่าอัตราความถูกต้องของการทำนายของโมเดล โดยในการวิจัยครั้งนี้ ผู้วิจัยตั้งเป้าหมายของค่าความถูกต้องไว้ที่ $> 90\%$ และจะทำการทดสอบค่าความถูกต้องในการทำนายของโมเดลที่เทรนด้วยวิธี Cross Validation โดยทำการแบ่งข้อมูลออกเป็น 2 ส่วน ได้แก่ ส่วนที่เอาไว้ใช้สำหรับการเทรนและอีกส่วนคือส่วนสำหรับการทดสอบ จะทำการสุ่มข้อมูลตามอัตราส่วนร้อยละ 90:10 และแบ่งข้อมูลสำหรับทำ K-Fold 5 Fold

3.4 การเปรียบเทียบประสิทธิภาพโมเดล

ในขั้นตอนการเปรียบเทียบประสิทธิภาพ ผู้วิจัยจะนำโมเดลที่ผ่านการเทรนทั้งหมด 3 โมเดล ได้แก่ LSTM, GRU, BiLSTM ซึ่งจะเปรียบเทียบประสิทธิภาพเรื่องของ ค่า Accuracy, ค่า Loss และ จำนวนรอบที่ใช้ในการเทรนโมเดล (epochs) เพื่อหาว่าโมเดลใด มีความแม่นยำมากที่สุด แล้วจะนำโมเดลที่มีความแม่นยำมากที่สุดนั้นมาทดสอบทำนายท่าทางภาษามือไทย

3.5 การทดสอบโมเดล

หลังจากได้รับโมเดลที่มีประสิทธิภาพที่ดีที่สุดแล้ว ผู้วิจัยจะนำโมเดลนั้นมาทดสอบด้วยวิดีโอที่เตรียมไว้สำหรับทดสอบ โดยประเภทของการทดสอบนั้นจะมีอยู่ 2 รูปแบบได้แก่ 1. เป็นคำศัพท์ 2. เป็นประโยค