

# **PROJECT - PART B: REPORT**

**CSI 4133 - COMP. METHODS IN PICTURE PROCESSING AND ANALYSIS**

**Fall 2022**

**School of Engineering and Computer Science**

**University of Ottawa**

**Course Coordinator: Dr. Jiying Zhao**

**Teaching Assistant: Christopher McIntyre Garciam**

Student Name: Jayden Lachhman

Student Number: 8791694

Submission date: 2022/12/02

## Introduction

This report will explore hand detection and tracking, finger detection and tracking, and neural network training with the overarching goal of designing a program capable of hand gesture recognition. This will be achieved by using a program that requires training data to learn properties of certain classifications of hand gestures so that it becomes better at classifying new data accurately.

## Approach

Given that hands and fingers come in different shapes, colors, and sizes, it would be impractical to depend on a comprehensive dataset to effectively analyze and correctly classify hand gestures. Instead, a machine learning approach (specifically a deep learning approach) will be used to teach our program how to learn hand gestures. This approach can be summarized by an iterative three-step process; the first involving the collection of training data which depends on the hand landmark model built into MediaPipe (a module within Python) to provide us with a framework from which we can extract information about the hands, the second involving the classification of this training data, and third the live analysis of a hand from the webcam of which comparisons to the training data are made to classify the new data accurately. The code references an existing repository for hand detection accessible here:

<https://github.com/kinivi/hand-gesture-recognition-mediapipe>

## Procedure

**Step 1:** The capture properties are initialized with reference to the device camera's width and height along with the tracking confidence threshold values which help us calibrate the MediaPipe model.

```
10 import cv2 as cv
11 import numpy as np
12 import mediapipe as mp
13
14 from utils import CvFpsCalc
15 from model import KeyPointClassifier
16 from model import PointHistoryClassifier
17
18
19 def get_args():
20     parser = argparse.ArgumentParser()
21
22     parser.add_argument("--device", type=int, default=0)
23     parser.add_argument("--width", help='cap width', type=int, default=960)
24     parser.add_argument("--height", help='cap height', type=int, default=540)
25
26     parser.add_argument('--use_static_image_mode', action='store_true')
27     parser.add_argument("--min_detection_confidence",
28                         help='min_detection_confidence',
29                         type=float,
30                         default=0.7)
31     parser.add_argument("--min_tracking_confidence",
32                         help='min_tracking_confidence',
33                         type=int,
34                         default=0.5)
35
36     args = parser.parse_args()
37
38     return args
```

Figure 1: Argument Initializations

```

41 def main():
42     # Argument parsing #####
43     args = get_args()
44
45     cap_device = args.device
46     cap_width = args.width
47     cap_height = args.height
48
49     use_static_image_mode = args.use_static_image_mode
50     min_detection_confidence = args.min_detection_confidence
51     min_tracking_confidence = args.min_tracking_confidence
52
53     use_brect = True
54
55     # Camera preparation #####
56     cap = cv.VideoCapture(cap_device)
57     cap.set(cv.CAP_PROP_FRAME_WIDTH, cap_width)
58     cap.set(cv.CAP_PROP_FRAME_HEIGHT, cap_height)
59
60     # Model load #####
61     mp_hands = mp.solutions.hands
62     hands = mp_hands.Hands(
63         static_image_mode=use_static_image_mode,
64         max_num_hands=2,
65         min_detection_confidence=min_detection_confidence,
66         min_tracking_confidence=min_tracking_confidence,
67     )

```

Figure 2: Model Initializations

### Hand Landmark Model

After the palm detection over the whole image our subsequent hand landmark [model](#) performs precise keypoint localization of 21 3D hand-knuckle coordinates inside the detected hand regions via regression, that is direct coordinate prediction. The model learns a consistent internal hand pose representation and is robust even to partially visible hands and self-occlusions.

To obtain ground truth data, we have manually annotated ~30K real-world images with 21 3D coordinates, as shown below (we take Z-value from image depth map, if it exists per corresponding coordinate). To better cover the possible hand poses and provide additional supervision on the nature of hand geometry, we also render a high-quality synthetic hand model over various backgrounds and map it to the corresponding 3D coordinates.

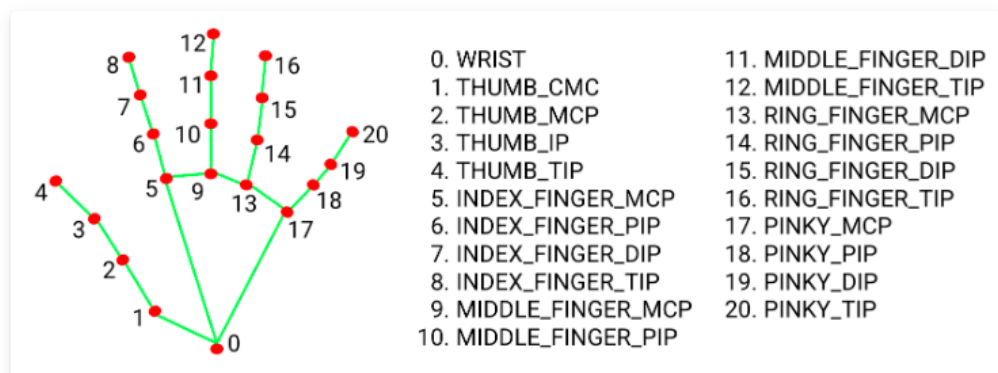


Figure 2: Hand Landmark Model

**Step 2:** In order for the program to classify the properties of input data, classes need to exist. Thus we create some .csv files which store the labels for each of the different classes the program will refer to when attempting to classify new hand gestures.

```

69     keypoint_classifier = KeyPointClassifier()
70
71     point_history_classifier = PointHistoryClassifier()
72
73     # Read labels #####
74     with open('model/keypoint_classifier/keypoint_classifier_label.csv',
75               encoding='utf-8-sig') as f:
76         keypoint_classifier_labels = csv.reader(f)
77         keypoint_classifier_labels = [
78             row[0] for row in keypoint_classifier_labels
79         ]
80     with open(
81         'model/point_history_classifier/point_history_classifier_label.csv',
82         encoding='utf-8-sig') as f:
83         point_history_classifier_labels = csv.reader(f)
84         point_history_classifier_labels = [
85             row[0] for row in point_history_classifier_labels
86         ]
87
88     # FPS Measurement #####
89     cvFpsCalc = CvFpsCalc(buffer_len=10)
90
91     # Coordinate history #####
92     history_length = 16
93     point_history = deque(maxlen=history_length)
94
95     # Finger gesture history #####
96     finger_gesture_history = deque(maxlen=history_length)

```

Figure 4: Keypoint Classifier Label

**Step 3:** Since MediaPipe requires the frames in RGB format and reads a mirrored reflection of the input video file, some presets are required to prepare the environment. A while loop is used since each frame needs to be sequentially acknowledged (meaning of *ret* variable), flipped, and converted from BGR to RGB.

```

99     mode = 0
100
101     while True:
102         fps = cvFpsCalc.get()
103
104         # Process Key (ESC: end) #####
105         key = cv.waitKey(10)
106         if key == 27: # ESC
107             break
108         number, mode = select_mode(key, mode)
109
110         # Camera capture #####
111         # ret is used to detect if a frame is available to read
112         ret, image = cap.read()
113         if not ret:
114             break
115         image = cv.flip(image, 1) # Mirror display
116         debug_image = copy.deepcopy(image)
117
118         # Detection implementation #####
119         image = cv.cvtColor(image, cv.COLOR_BGR2RGB)
120
121         image.flags.writeable = False
122         results = hands.process(image)
123         image.flags.writeable = True

```

Figure 5: Environment Initialization

**Step 4:** MediaPipe and the Hand Landmark Model help us generate numerical information about the input data, but this information needs to be reformatted in order to handle it properly. This starts with the *hand\_landmarks* list which reads the x, y, and z coordinates of the detected hand in each frame. These coordinates do not signify a hand gesture independently, but in relation to one another they can. Thus, after creating the region of interest (ROI) referred to by the *brect* variable, the raw pixel coordinates are converted into relative coordinates where eventually the distance between points can be compared to distances between points from the training data in accordance with the Hand Landmark Model. These pre-processed points are then entered into the .csv.

```

125 # hand_landmarks represents the x, y, z coordinates of the detected hand in each frame
126 # handedness represents the orientation of a hand to track it under movement
127 if results.multi_hand_landmarks is not None:
128     for hand_landmarks, handedness in zip(results.multi_hand_landmarks,
129                                           results.multi_handedness):
130         # The following print statement prints the x, y, z coordinates of the detected hand in each frame
131         # print(hand_landmarks)
132
133         # Bounding box calculation
134         brect = calc_bounding_rect(debug_image, hand_landmarks)
135         # Landmark calculation
136         # landmark_list represents the list of hand_landmarks corresponding to the 21 landmarks on the hand landmark model
137         landmark_list = calc_landmark_list(debug_image, hand_landmarks)
138         # The following print statement prints the xy-coordinates corresponding to the wrist landmark in the hand landmark model
139         # print(landmark_list[0])
140
141         # Conversion to relative coordinates / normalized coordinates
142         # This conversion is important because without normalizing the coordinates, the program would be reading exact pixel
143         # coordinates instead of interpreting them relative to each other. The distance between landmark points characterizes
144         # gestures most distinctly, otherwise, a hand on one side of the screen might be recognized differently to that same
145         # hand doing the same gesture on the other side of the screen
146         pre_processed_landmark_list = pre_process_landmark(
147             landmark_list)
148         pre_processed_point_history_list = pre_process_point_history(
149             debug_image, point_history)
150         # Write to the dataset file
151         logging_csv(number, mode, pre_processed_landmark_list,
152                    pre_processed_point_history_list)

```

Figure 6: Input Pre-Processing

**Step 5:** Extended from the previous step, pre-processing the coordinates helps us more accurately classify the hand gesture they might be demonstrating. This can be done by assigning a *base* variable to the wrist point on the Hand Landmark Model which serves to be the point the other points refer to. This is a crucial step because the raw coordinates present a large spectrum of potential interpretation, which is proportionally error prone. In this sense, if one were to attempt to classify a gesture solely based on the coordinates of each point independently, hand movement across the camera could produce false-positives/true-negatives. Using the wrist as a base means we measure movement from a relatively static point, which is less error prone.

```

245 def pre_process_landmark(landmark_list):
246     temp_landmark_list = copy.deepcopy(landmark_list)
247
248     # Convert to relative coordinates
249     # This sets the landmark point corresponding to the wrist point to zero (0) which is our base value. This is important
250     # because now the distance of the other points on the hand landmark model are calculated relative to this wrist point
251     # so that our gestures are detectable regardless of where the hand is on the screen, which is more effective than trying
252     # to detect gestures using independent pixel values
253     base_x, base_y = 0, 0
254     for index, landmark_point in enumerate(temp_landmark_list):
255         if index == 0:
256             base_x, base_y = landmark_point[0], landmark_point[1]
257
258             temp_landmark_list[index][0] = temp_landmark_list[index][0] - base_x
259             temp_landmark_list[index][1] = temp_landmark_list[index][1] - base_y

```

Figure 7: Conversion to Wrist Relative Coordinates

**Step 6:** Extending again from the previous step, the relative coordinates are more useful than in their previous format, but can be improved still. The landmark displacement values can be made more distinctive by normalizing each to a range between -1 and 1 instead of calculating each across arbitrary ranges, and after this, the pre-processing stage is complete.

```

261 # Convert to a one-dimensional list
262 temp_landmark_list = list(
263     itertools.chain.from_iterable(temp_landmark_list))
264
265 # Normalization
266 # The list of landmark values are relative to the wrist point, but the format of these values can be made easier
267 # still for the program to recognize hand gestures from by normalizing the values to numbers between -1 and 1
268 # by dividing their the values by the max in the list. This helps the program interpret more distinctly the points
269 # of the landmark model when classifying the gestures.
270 max_value = max(list(map(abs, temp_landmark_list)))
271
272 def normalize_(n):
273     return n / max_value
274
275 temp_landmark_list = list(map(normalize_, temp_landmark_list))
276
277 return temp_landmark_list

```

Figure 8: Normalization

**Step 7:** Each of the distinct hand gestures refers to an index in the *keypoint\_classifier\_label.csv* file from earlier. This is where they are matched and the resulting hand gesture is returned.

```

154 # Hand sign classification
155 # This signifies that the neural network responds with the classification of "pointer" for the index of 2, which corresponds
156 # to the index of each gesture classification in the other file
157 hand_sign_id = keypoint_classifier(pre_processed_landmark_list)
158 if hand_sign_id == 2: # One gesture
159     point_history.append(landmark_list[8])
160 else:
161     point_history.append([0, 0])
162
163 # Finger gesture classification
164 finger_gesture_id = 0
165 point_history_len = len(pre_processed_point_history_list)
166 if point_history_len == (history_length * 2):
167     finger_gesture_id = point_history_classifier(
168         pre_processed_point_history_list)
169
170 # Calculates the gesture IDs in the latest detection
171 finger_gesture_history.append(finger_gesture_id)
172 most_common_fg_id = Counter(
173     finger_gesture_history).most_common()
174
175 # Drawing part
176 debug_image = draw_bounding_rect(use_brect, debug_image, brect)
177 debug_image = draw_landmarks(debug_image, landmark_list)
178 debug_image = draw_info_text(
179     debug_image,
180     brect,
181     handedness,
182     keypoint_classifier_labels[hand_sign_id],
183     point_history_classifier_labels[most_common_fg_id[0][0]],

```

Figure 9: Hand and Finger Classification

```

model > keypoint_classifier > keypoint_classifier_label.csv
1   Five
2   Zero
3   One
4   Two
5   Three
6   Four

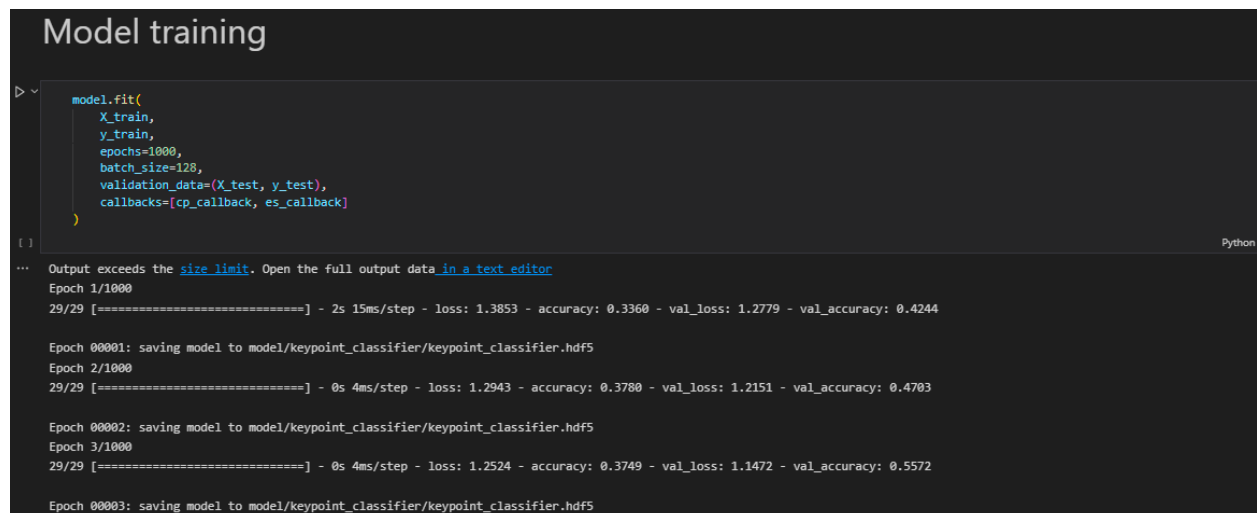
```

Figure 10: Classifier Lable CSV

**Step 8:** Collecting training data is important for supplementing our neural network with more data to use in its pool to invoke when classifying a hand gesture. So instead of importing a dataset, this program allows the efficient creation of a dataset. Below is a figure representing different modes once the camera has begun capturing frames, and inputting 'k' will allow the user to log the hand landmarks currently in frame. Then by pressing the number corresponding to the index of the class the user would like to populate training data for, the coordinates of the landmarks will be logged into the .csv file accessible through the Jupyter notebook where similar data for other classes are stored.

```
198 def select_mode(key, mode):
199     number = -1
200     if 48 <= key <= 57: # 0 ~ 9
201         number = key - 48
202     if key == 110: # n
203         mode = 0
204     if key == 107: # k
205         mode = 1
206     if key == 104: # h
207         mode = 2
208     return number, mode
```

Figure 9: Capture Mode

The image shows a Jupyter Notebook interface with a dark theme. At the top, the title 'Model training' is displayed. Below the title, a code cell contains a `model.fit()` call with parameters: `X_train`, `y_train`, `epochs=1000`, `batch_size=128`, `validation_data=(X_test, y_test)`, and `callbacks=[cp_callback, es_callback]`. The output of the cell shows the training progress over three epochs. Each epoch summary includes the number of steps (29/29), time per step (15ms, 4ms, 4ms), loss, accuracy, validation loss, and validation accuracy. The model is saved to `model/keypoint_classifier/keypoint_classifier.hdf5` at the end of each epoch. A message at the top of the output indicates that the output exceeds the size limit and suggests opening the full output data in a text editor.

```
model.fit(
    X_train,
    y_train,
    epochs=1000,
    batch_size=128,
    validation_data=(X_test, y_test),
    callbacks=[cp_callback, es_callback]
)
```

Output exceeds the [size limit](#). Open the full output data [in a text editor](#)

Epoch 1/1000  
29/29 [=====] - 2s 15ms/step - loss: 1.3853 - accuracy: 0.3360 - val\_loss: 1.2779 - val\_accuracy: 0.4244

Epoch 00001: saving model to model/keypoint\_classifier/keypoint\_classifier.hdf5  
Epoch 2/1000  
29/29 [=====] - 0s 4ms/step - loss: 1.2943 - accuracy: 0.3780 - val\_loss: 1.2151 - val\_accuracy: 0.4703

Epoch 00002: saving model to model/keypoint\_classifier/keypoint\_classifier.hdf5  
Epoch 3/1000  
29/29 [=====] - 0s 4ms/step - loss: 1.2524 - accuracy: 0.3749 - val\_loss: 1.1472 - val\_accuracy: 0.5572

Epoch 00003: saving model to model/keypoint\_classifier/keypoint\_classifier.hdf5

Figure 10: Model Training in Jupyter Notebook

The seven (7) steps detailed above describe the most functionally/conceptually critical points of the program but there are a number of helper functions that are instrumental in terms of coherence. These will be summarized below but separated from the linear sequence of principal steps detailed above.

The following code snippet (Figure 11) refers to the *logging\_csv()* function where the training data for the hand gesture classification is stored. The more data entered, the more it helps increase the accuracy level of our gesture recognition program.

```
303 def logging_csv(number, mode, landmark_list, point_history_list):
304     if mode == 0:
305         pass
306     if mode == 1 and (0 <= number <= 9):
307         csv_path = 'model/keypoint_classifier/keypoint.csv'
308         with open(csv_path, 'a', newline='') as f:
309             writer = csv.writer(f)
310             writer.writerow([number, *landmark_list])
311     if mode == 2 and (0 <= number <= 9):
312         csv_path = 'model/point_history_classifier/point_history.csv'
313         with open(csv_path, 'a', newline='') as f:
314             writer = csv.writer(f)
315             writer.writerow([number, *point_history_list])
316     return
317
```

Figure 11: Logging CSV

The following code snippet (Figure 12) refers to the *draw\_landmarks()* function which uses the Hand Landmark Model referenced earlier to identify the landmarks composing each finger.

```
319 def draw_landmarks(image, landmark_point):
320     if len(landmark_point) > 0:
321         # Thumb
322         cv.line(image, tuple(landmark_point[2]), tuple(landmark_point[3]),
323                 (0, 0, 0), 6)
324         cv.line(image, tuple(landmark_point[2]), tuple(landmark_point[3]),
325                 (255, 255, 255), 2)
326         cv.line(image, tuple(landmark_point[3]), tuple(landmark_point[4]),
327                 (0, 0, 0), 6)
328         cv.line(image, tuple(landmark_point[3]), tuple(landmark_point[4]),
329                 (255, 255, 255), 2)
330
331         # Index finger
332         cv.line(image, tuple(landmark_point[5]), tuple(landmark_point[6]),
333                 (0, 0, 0), 6)
334         cv.line(image, tuple(landmark_point[5]), tuple(landmark_point[6]),
335                 (255, 255, 255), 2)
336         cv.line(image, tuple(landmark_point[6]), tuple(landmark_point[7]),
337                 (0, 0, 0), 6)
338         cv.line(image, tuple(landmark_point[6]), tuple(landmark_point[7]),
339                 (255, 255, 255), 2)
340         cv.line(image, tuple(landmark_point[7]), tuple(landmark_point[8]),
341                 (0, 0, 0), 6)
342         cv.line(image, tuple(landmark_point[7]), tuple(landmark_point[8]),
343                 (255, 255, 255), 2)

```

Figure 12: Landmark Detection



The following code snippet (Figure 13) refers to the *draw\_bounding\_rect()* function which calculates the rectangle bounding the user's hand. The camera's dimensions include the user's hand, but there is excess information stored there as well, and thus the program needs to focus on extracting only pertinent data. This pertinence can be defined using a rectangular bound around the user's hand.

```

507 def draw_bounding_rect(use_brect, image, brect):
508     if use_brect:
509         # Outer rectangle
510         cv.rectangle(image, (brect[0], brect[1]), (brect[2], brect[3]),
511                       (0, 0, 0), 1)
512
513     return image
514
515
516 def draw_info_text(image, brect, handedness, hand_sign_text,
517                   finger_gesture_text):
518     cv.rectangle(image, (brect[0], brect[1]), (brect[2], brect[1] - 22),
519                 (0, 0, 0), -1)
520
521     info_text = handedness.classification[0].label[0:]
522     if hand_sign_text != "":
523         info_text = info_text + ':' + hand_sign_text
524     cv.putText(image, info_text, (brect[0] + 5, brect[1] - 4),
525               cv.FONT_HERSHEY_SIMPLEX, 0.6, (255, 255, 255), 1, cv.LINE_AA)
526
527     if finger_gesture_text != "":
528         cv.putText(image, "Finger Gesture:" + finger_gesture_text, (10, 60),
529                   cv.FONT_HERSHEY_SIMPLEX, 1.0, (0, 0, 0), 4, cv.LINE_AA)
530         cv.putText(image, "Finger Gesture:" + finger_gesture_text, (10, 60),
531                   cv.FONT_HERSHEY_SIMPLEX, 1.0, (255, 255, 255), 2,
532                   cv.LINE_AA)
533
534     return image

```

Figure 13: Region of Interest

The following code snippet (Figure 14) refers to a couple of *draw...()* functions. Each adds a different textual element to the frame when the camera is recording, that element related to the measured frame-rate, mode of capture, etc.

```

537 def draw_point_history(image, point_history):
538     for index, point in enumerate(point_history):
539         if point[0] != 0 and point[1] != 0:
540             cv.circle(image, (point[0], point[1]), 1 + int(index / 2),
541                       (152, 251, 152), 2)
542
543     return image
544
545
546 def draw_info(image, fps, mode, number):
547     cv.putText(image, "FPS:" + str(fps), (10, 30), cv.FONT_HERSHEY_SIMPLEX,
548               1.0, (0, 0, 0), 4, cv.LINE_AA)
549     cv.putText(image, "FPS:" + str(fps), (10, 30), cv.FONT_HERSHEY_SIMPLEX,
550               1.0, (255, 255, 255), 2, cv.LINE_AA)
551
552     mode_string = ['Logging Key Point', 'Logging Point History']
553     if 1 <= mode <= 2:
554         cv.putText(image, "MODE:" + mode_string[mode - 1], (10, 90),
555                   cv.FONT_HERSHEY_SIMPLEX, 0.6, (255, 255, 255), 1,
556                   cv.LINE_AA)
557         if 0 <= number <= 9:
558             cv.putText(image, "NUM:" + str(number), (10, 110),
559                       cv.FONT_HERSHEY_SIMPLEX, 0.6, (255, 255, 255), 1,
560                       cv.LINE_AA)
561     return image
562
563
564 if __name__ == '__main__':
565     main()
566

```

Figure 14: Info

The following snippet from the jupyter file shows how the new data inputs are being synthesized with the existing training data to increase the accuracy of this particular hand gesture classification.

```
Epoch 58/1000
21/30 [=====>.....] - ETA: 0s - loss: 0.6496 - accuracy: 0.7612
Epoch 58: saving model to model/keypoint_classifier/keypoint_classifier.hdf5
30/30 [=====] - 0s 5ms/step - loss: 0.6706 - accuracy: 0.7511 - val_loss: 0.3491 - val_accuracy: 0.8909
Epoch 59/1000
20/30 [=====>.....] - ETA: 0s - loss: 0.6595 - accuracy: 0.7523
Epoch 59: saving model to model/keypoint_classifier/keypoint_classifier.hdf5
30/30 [=====] - 0s 5ms/step - loss: 0.6592 - accuracy: 0.7542 - val_loss: 0.3456 - val_accuracy: 0.8973
Epoch 60/1000
19/30 [=====>.....] - ETA: 0s - loss: 0.6836 - accuracy: 0.7405
Epoch 60: saving model to model/keypoint_classifier/keypoint_classifier.hdf5
30/30 [=====] - 0s 5ms/step - loss: 0.6667 - accuracy: 0.7495 - val_loss: 0.3440 - val_accuracy: 0.8973
Epoch 61/1000
22/30 [=====>.....] - ETA: 0s - loss: 0.6711 - accuracy: 0.7454
Epoch 61: saving model to model/keypoint_classifier/keypoint_classifier.hdf5
30/30 [=====] - 0s 6ms/step - loss: 0.6724 - accuracy: 0.7402 - val_loss: 0.3328 - val_accuracy: 0.9021
```

Figure 15: Learning

## Data

The following three (3) images correspond to the hand gesture classifications determined by the program. Each hand gesture refers to a certain number of fingers raised.

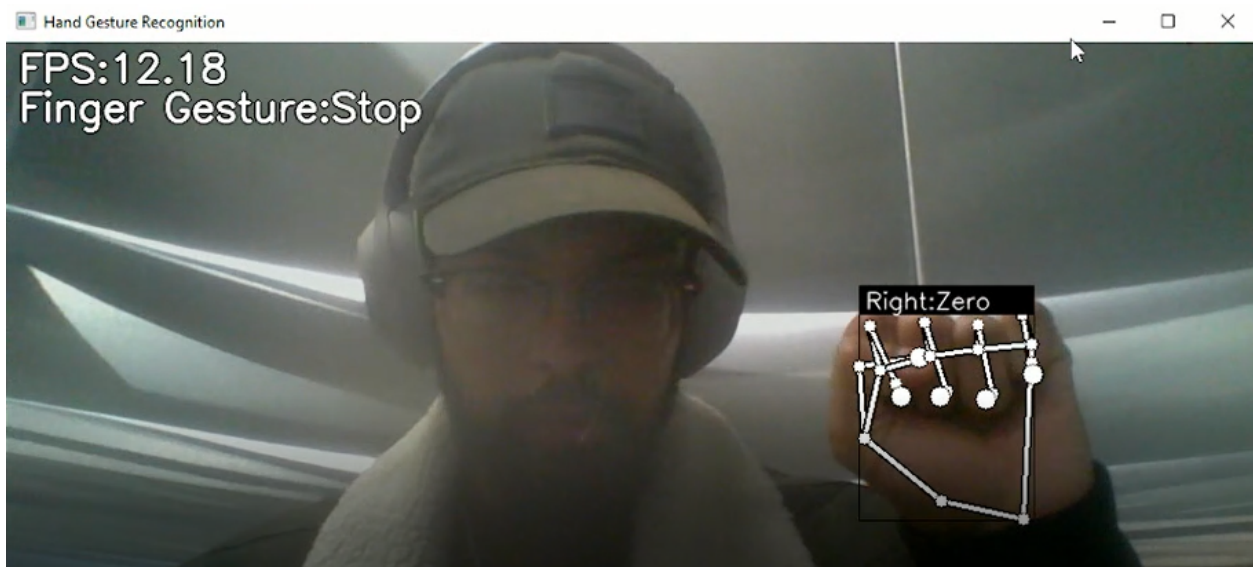


Figure 1: Zero Fingers Raised

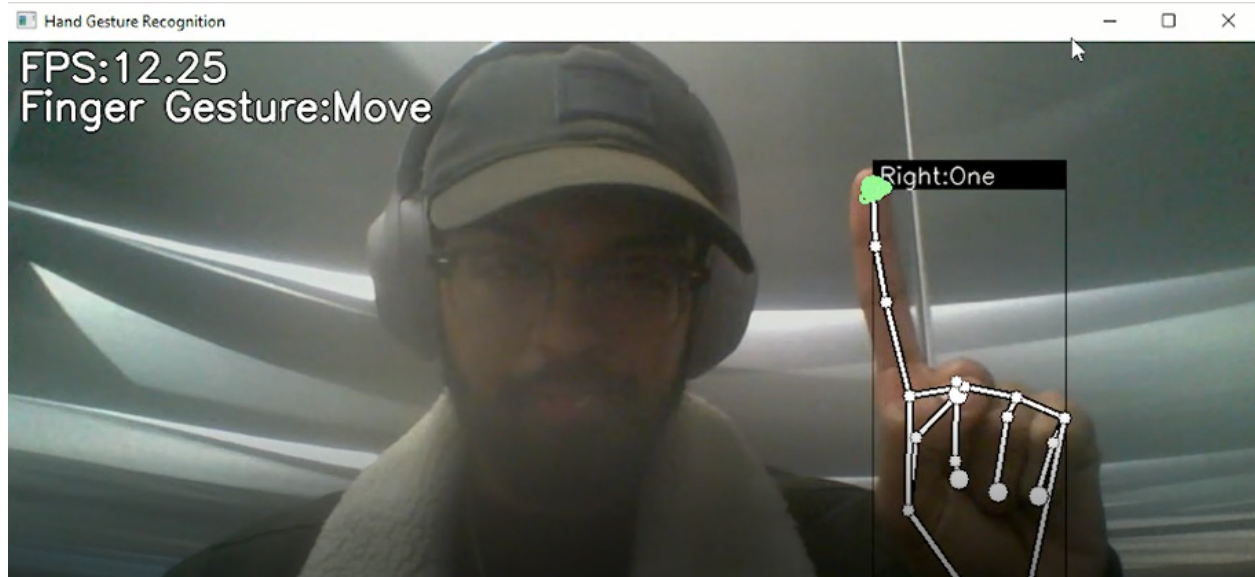


Figure 2: One Finger Raised

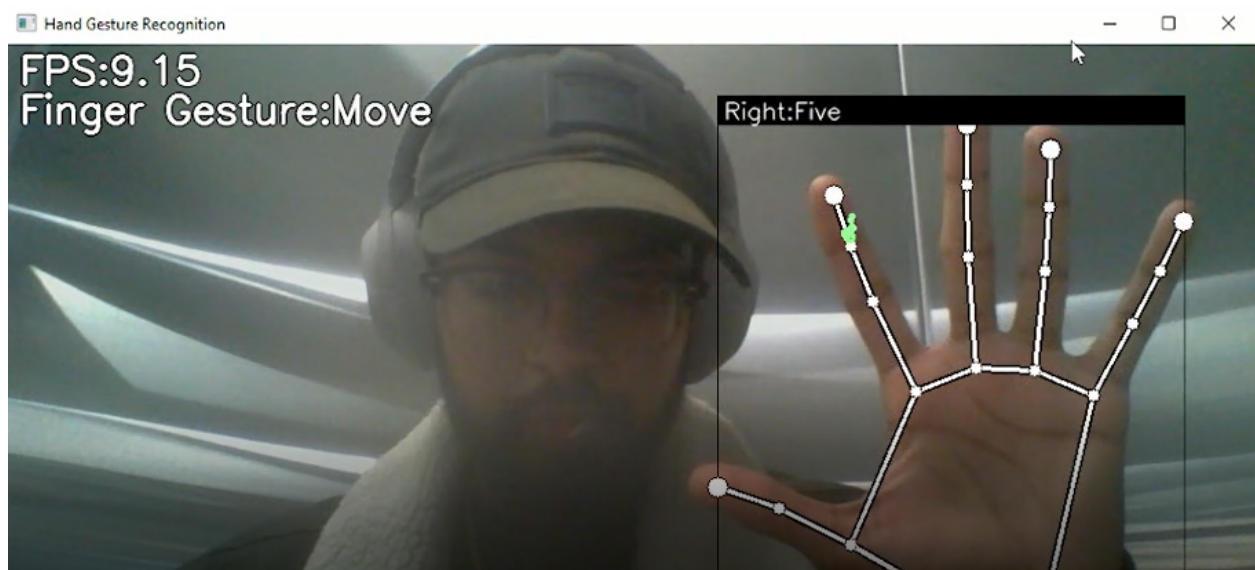


Figure 3: Five Fingers Raised

## Conclusion

Using the MediaPipe module with reference to the Hand Landmark Model, a couple of .csv files for logging and storing training data for a neural network, and a jupyter notebook for synchronizing new data inputs (frames) with others of the same hand gesture classification, we were able to design a program that helps us analyze hand gestures from a live video feed. The accuracy of