

Behavior recognition using Kinect

Jiangjian Guo

Qiming Liu

Weihao Yan

January 10, 2019

Abstract

Human action recognition has great practical value and is currently a hot research field of artificial intelligence and machine learning. Our project intends to use the appropriate method to enable the machine to recognize the human body in real time through Kinect and mark the corresponding action.

We first perform feature extraction and preprocessing on the dataset through Openpose, which converts the video samples into coordinates of key points of the human body. We then used a convolutional neural network approach to make the machine learn the characteristics of the relevant actions and achieved good results in real-time testing. In order to overcome the problems of insufficient data and recognition errors, we also built our own data sets and built a new deep convolutional network. By comparison, the test results have been greatly improved.

To improve the practical value of the project, we added the function of multi-person action recognition in real-time action recognition, and intend to add depth pictures for training in subsequent work to improve the recognition accuracy.

In this paper, we will detail the complete process and implementation of the project, and give the test results. Of course, there are still a lot of things to improve in the project, we will also introduce them and try to propose possible solutions.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction and Motivation | 3 |
| 2 | Background | 3 |
| 2.1 | Prior Work | 3 |
| 2.1.1 | Action Representation | 3 |
| 2.1.2 | Action Classification | 4 |
| 2.1.3 | Deep Network | 4 |
| 2.2 | Research Challenges | 4 |
| 2.2.1 | Intra- and Inter-class Variations | 4 |
| 2.2.2 | Cluttered Background and Camera Motion | 4 |
| 2.2.3 | Insufficient Annotated Data | 4 |
| 3 | Implementation Details | 4 |
| 3.1 | Dataset | 4 |
| 3.2 | Feature Extraction Using Openpose | 5 |
| 3.3 | Introduction to Deep Learning Network | 6 |
| 3.4 | New Dataset: Action-209 | 6 |
| 3.5 | New Deep Network Based on Actions-209 | 9 |
| 4 | Experiments | 10 |
| 4.1 | Real-time Video Sampling | 10 |
| 4.2 | Parameters Tunning | 10 |
| 4.2.1 | The Tunning Process and Experiment Result of First Net | 10 |
| 4.2.2 | The Tunning Process and Experiment Result of Second Net | 11 |
| 5 | Source code | 13 |
| 5.1 | Dataset Preprocessing Using Openpose | 13 |
| 5.2 | Training Process | 15 |
| 5.3 | Create Dataset Action-209 | 20 |
| 5.4 | Real-time Recognition | 27 |
| 6 | Discussion and Conclusions | 34 |
| 6.1 | Discussion | 34 |
| 6.1.1 | Missing Data Processing | 34 |
| 6.1.2 | The Use of Depth Information | 34 |
| 6.1.3 | Our New Dataset | 34 |
| 6.1.4 | Real-time Video Sampling | 36 |
| 6.2 | Conclusion | 36 |

1 Introduction and Motivation

Human behavior recognition is a bench of computer version, which has been a popular research area in fields of artificial intelligence and machine learning in recent years. It has so many applications such as human-computer interfaces, content-based video retrieval and public security guard etc. At the same time, behavior recognition is a very challenging task because given human action can have different meanings. Different people may also perform an action in dynamic ways. Moreover, changing environment element such as light, background, occlusion and view direction may bring much uncertainty to the recognizing processes.

Let's imagine such a situation happening in the coming future. Games will no longer be limited to computer and console games. Motion sensing games will become more popular concerning health and interaction requirement. However, good motion sensing games need an action recognition system to recognize different actions of player. It should be real-time and precise so that everyone can enjoy it. Human action recognition is also a significant part in AI assistant, which will become more and more popular in next few years. Evenmore, the recognizing system can be applied on many domains. For example, we can monitor old people's action to avoid some emergent accidents. We can also apply it on robots so that they can interact with human better. A well-performed action recognition system can profoundly change our life. Inspired by the amazing usage of action recognition, our team was thinking about building a real-time human behavior recognition system.

Our project is aimed at making computer be aware of human action and recognize what behaviors the human is performing in an effectively and efficiently way. We are required to use Microsoft Kinect, which is a series of peripheral containing RGB camera, depth camera, microphones, etc. The final project result is supposed to output the bounding box of human and related action labels in a real-time video.

In our project, we used Openpose to preprocess and do feature extration to our original dataset. Many algorithms including C3D, C2D, light flow have been tried. We also created our own dataset to have a better result of recognition. Detailed information and result of our project will be introduced in this report. Also, the follow-up research plan and some defect will also be introduced.

2 Background

2.1 Prior Work

2.1.1 Action Representation

Holistic representation methods capture the motion information of the entire human subject. Local representations only identify local regions having salient motion information.

2.1.2 Action Classification

Action classifiers can be roughly divided into the following categories: Classification, Sequential Approaches, Space-time Approaches, Partbased Approaches, Manifold Learning Approaches, Mid-Level Feature Approaches and Feature Fusion Approaches

2.1.3 Deep Network

In recent years, feature learning using deep learning techniques has been receiving increasing attention due to their ability of designing powerful features that can be generalized very well [2, 3, 5] . Recent deep networks [2, 4, 6] have achieved surprisingly high recognition performance on a variety of action datasets.

2.2 Research Challenges

Despite significant progress has been made in human action recognition and prediction, state-of-the-art algorithms still misclassify actions due to several major challenges.

2.2.1 Intra- and Inter-class Variations

People behave differently for the same actions. Videos in the same action can be captured from various viewpoints

2.2.2 Cluttered Background and Camera Motion

Most of existing activity features such as histograms of oriented gradient and interest points also encode background noise, and thus degrade the recognition performance. Camera motion is another factor that should be considered in real-world applications. Due to significant camera motion, action features cannot be accurately extracted.

2.2.3 Insufficient Annotated Data

Even though existing action recognition approaches have shown impressive performance on small-scale datasets in laboratory settings, it is really challenging to generalize them to real-world applications due to their inability of training on largescale datasets.

3 Implementation Details

3.1 Dataset

Finding a suitable dataset for our algorithm is a time-consuming task. Most datasets available on the internet are either too small or without acceptable and correct labels. We spent several days and finally decided to use FLORENCE 3D ACTIONS DATASET as our preliminary dataset. This dataset includes 9 actions: wave, drink from a bottle, answer phone,

clap, tight lace, sit down, stand up, read watch, bow. Each action is performed by ten people and each person should perform an action for 2 or 3 times. Some screenshots are pictured in Figure 1.



Figure 1: Some screenshots of dataset.

Obviously, this dataset is very small but also well-labeled. We decide to use it to confirm our algorithm first, and then replace the dataset with a much bigger one. We created our own dataset during project time, related work will be introduced later in this report.

3.2 Feature Extraction Using Openpose

In FLORENCE 3D ACTIONS DATASET, each sample is in video form, each video has 10 to 25 frames, each frame has 600×480 pixels, and each pixel has three channels (RGB). The total dimension of one sample is extremely large, we cannot get a good result in just 215 samples, which is the main challenge at the beginning of our project. Clearly, we need dimensionality reduction, feature extraction and data preprocessing.

We decided to use Openpose [7, 8, 9, 10] as a tool to help us with this essential work. Openpose represents a real-time multi-person system to detect human body, hand, facial, and foot key points on single images. It used deep net method (VGG, Resnet, Mobilenet, etc.) and Part Affinity Fields (PAFs) to estimate one's body key point. To learn more about openpose, you can visit its website on Github.

Note that we only need to load the pretrained model in Openpose to estimate key point of body in a picture. Openpose offered several models which differs in the recognizing speed, recognizing accuracy and the amount of key point. To fit our project task, we decide to use the MPI model. This model will output 15 key points of one's body. The test output is showed in Figure 2. When the model cannot recognize key parts of body, it will return a None type variable, which represents the part is missing.

Clearly seen from the test result, Openpose can transfer a high dimension picture to a series of key point positions, which in very low dimension. Furthermore, Openpose also helps us to extract features from the original graph, so that we can represent actions the human is playing with some point positions. Learning the moving law of key points is also learning how the actual human behavior looks like.

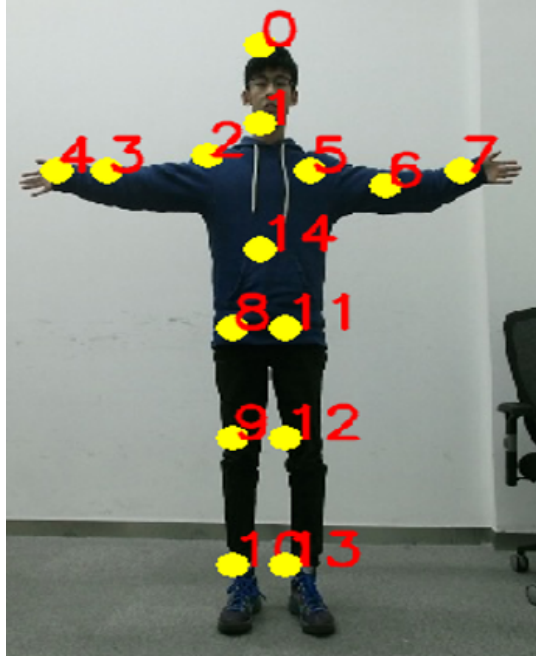


Figure 2: Some screenshots of dataset.

In practice, we divide a sample video into 8 frames, each frame can have X and Y position of 15 key points using Openpose. We spliced the position of key points in the order of frame occurred. In this way we got the preprocess new dataset.

3.3 Introduction to Deep Learning Network

Considering our data's shape, there are 8 frames in one action and in each frame there are 15 key points which have x and y two dimensions. Aftering reading this reference[1], we decide to use 2D Convolution. Specifically, our first net is shown in Figure 3.

In the net, x and y are two different channels while 8 frames and 15 key points are two dimensions to be convoluted. The convolution structure is shown in Figure 4.

Running the preprocessed dataset on this net, we finally got around 80% accuracy, detailed training process and result will be introduce at the Experiment part.

3.4 New Dataset: Action-209

In practice, we found that the model is not good enough because we still get some errors in real-time recognition. We tuned the parameters to improve the result but the test accuracy is always around 80%, which means the model need better generalization ability. We believe it is because the datasets we use is too small and monotonous. So we decided to build our own datasets. We call it Action-209 since it is collected in Lizhengdao library room 209.

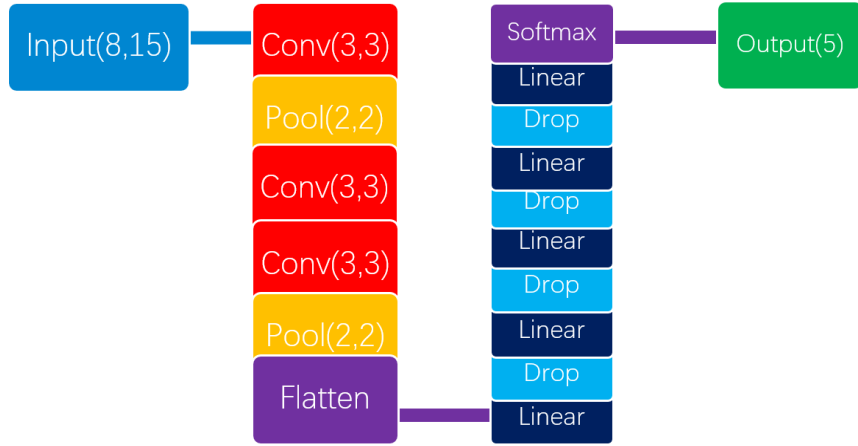


Figure 3: The structure of our first net

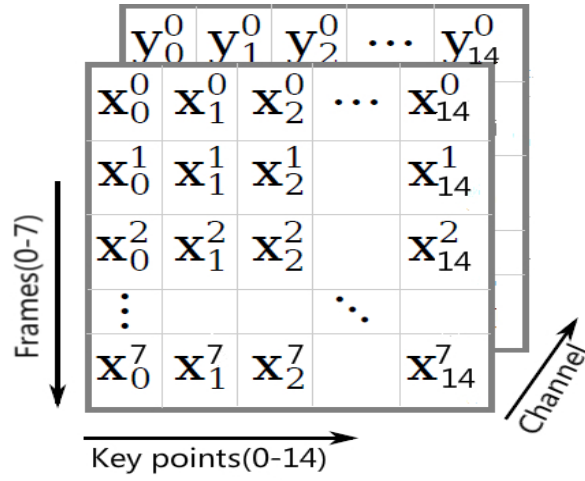


Figure 4: The convolution structure

We use the function of Pykinect API to catch real-time RGB and depth frames. Sampling frequency is fixed so that trained model will have a stable performance.

We designed the dataset according to the shortcomings of the previous one.

| FLORENCE 3D ACTIONS DATASET | Action-209 |
|-----------------------------|---------------------------------|
| Only 5 high-quality actions | Expanded to 9 actions |
| 25 series for an action | Nearly 240 series of images |
| Always from starts to ends | Might start at any time |
| Only has 8 frames | Collect 16 frames in a row |
| Low frequency | More detailed information |
| Always act in the center | Collect images of all positions |
| Similar speed when acting | Different speed when acting |

Table 1: The tuning process of parameters

We collect 240 samples for one action, each of them has 16 frames. It is much larger and better than the previous one and we did get a more precise model trained with it.



Figure 5: Action-209 dataset created by our team.

The datasets are collected in a tiny room containing data of only 3 group members so it still need to be improved. We will talk about it in Discussion part.

3.5 New Deep Network Based on Actions-209

Owing to more data, more frames and more action types, we need to update our first net to a more complex one. After adjusting the structure, the new net is shown in Figure 6.

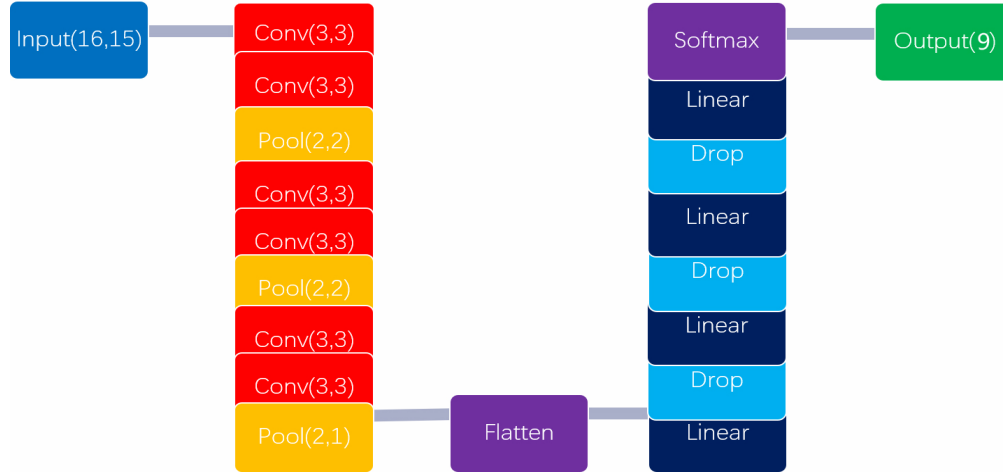


Figure 6: The structure of our second net

The convolution structure is updated to picture shown in Figure 7.

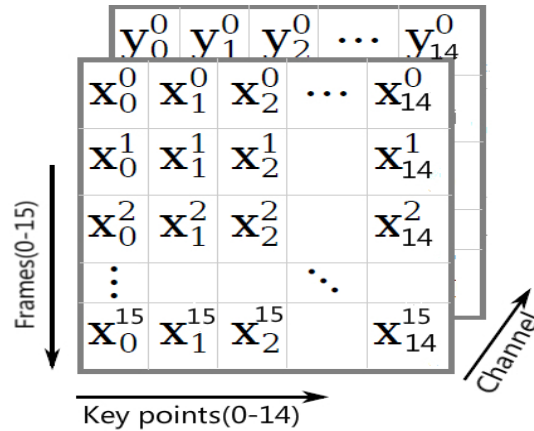


Figure 7: The convolution structure

Running the preprocessed dataset on this net, we finally got around 95% accuracy, detailed training process and result will be introduced at the Experiment part.

4 Experiments

4.1 Real-time Video Sampling

We planned to put the model into practical applications, which means we want to do real-time recognition with frames from Kinect rather than a preprocessed video. Fortunately, we found an API called pykinect2, which can catch RGB-D frame and joints coordinates instantly. So, while running the demo, we can have real-time information of the human in front of Kinect.

In real-time application, we keep pace with the sampling of dataset to get instant coordinates. We collect information of 16 frames and preprocess it to fit the model. The network then will give us a prediction label which is the action label we will update with.

Basically, we got a pretty good recognition system which can do well in real-time running. It can accurately distinguish 9 different actions right now.

What's more, we make our idea of two-people action recognition come true. When there are two people in front of Kinect doing actions, our algorithms can recognize them respectively. In general, we can increase the number of people to 6 at most.

But it still makes some mistakes during testing. We will talk about it in Discussion.

4.2 Parameters Tunning

4.2.1 The Tunning Process and Experiment Result of First Net

The main tuning parameters of our net are max_iterations, batch_size and learning_rate. We tried many combinations of parameters, which has been listed in Table 2.

| max_iterations | batch_size | learning_rate | train_accuracy(%) | test_accuracy(%) |
|----------------|------------|---------------|-------------------|------------------|
| 50 | 4 | 0.001 | 64.00 | 69.47 |
| 100 | 4 | 0.001 | 82.67 | 76.84 |
| 150 | 4 | 0.0005 | 95.20 | 74.73 |
| 150 | 5 | 0.0005 | 96.53 | 81.05 |
| 150 | 10 | 0.0005 | 93.87 | 84.21 |
| 200 | 10 | 0.0005 | 72.53 | 66.31 |
| 200 | 10 | 0.0001 | 66.13 | 52.63 |
| 200 | 20 | 0.0005 | 88.53 | 70.52 |
| 200 | 20 | 0.0001 | 59.46 | 46.32 |
| 200 | 15 | 0.0005 | 95.20 | 87.37 |

Table 2: The tuning process of parameters

According to Table 2, we chose max_iterations = 200, batch_size = 15 and learning_rate = 0.0005, the training process is shown in Figure 8.

From the result we know that after 200 iterations, the net's training accuracy is above 95% and testing accuracy is more than 85%, which seems a good result. However, due to

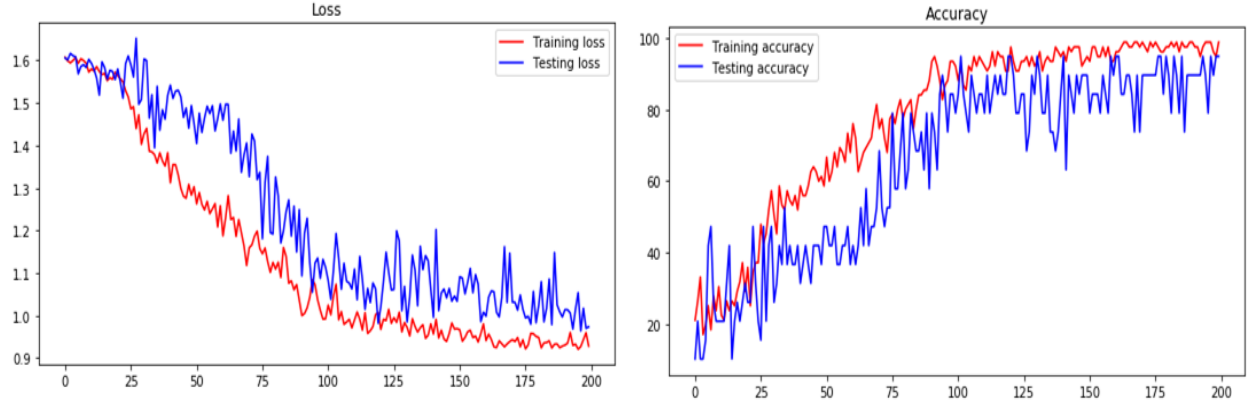


Figure 8: The variation trend of loss and accuracy of final version in the training process

small size of dataset and randomness of the training process, the result of training varies from time to time. Sometimes the testing accuracy is below 75%, which is not good enough.



Figure 9: The accuracy varies from time to time

4.2.2 The Tuning Process and Experiment Result of Second Net

We tried many combinations of parameters, which has been listed in Table 3.

Finally we chose $\text{max_iterations} = 300$, $\text{batch_size} = 4$ and $\text{learning_rate} = 0.00005$, the training process is shown in Figure 10.

From the result we can see that after 300 iterations, both the training accuracy and testing accuracy of our second net are higher than 90% which shows good performance and the curves converge more stably. What's more, after running the training process for several times, the accuracies are always higher than 90%.

We also tried to use RandomForest to get the prediction model. After tuning the hyper-parameters in RandomForest Classifier, the final result is shown in Figure 11:

| max_iterations | batch_size | learning_rate | train_accuracy(%) | test_accuracy(%) |
|----------------|------------|---------------|-------------------|------------------|
| 100 | 4 | 0.001 | 15.65 | 15.64 |
| 100 | 4 | 0.0005 | 17.47 | 17.04 |
| 150 | 4 | 0.0005 | 76.66 | 86.59 |
| 150 | 4 | 0.0001 | 92.73 | 93.30 |
| 200 | 10 | 0.0001 | 94.56 | 93.97 |
| 300 | 10 | 0.00005 | 86.44 | 85.47 |
| 300 | 10 | 0.00001 | 65.13 | 67.60 |
| 300 | 4 | 0.00005 | 96.51 | 95.81 |

Table 3: The tuning process of parameters

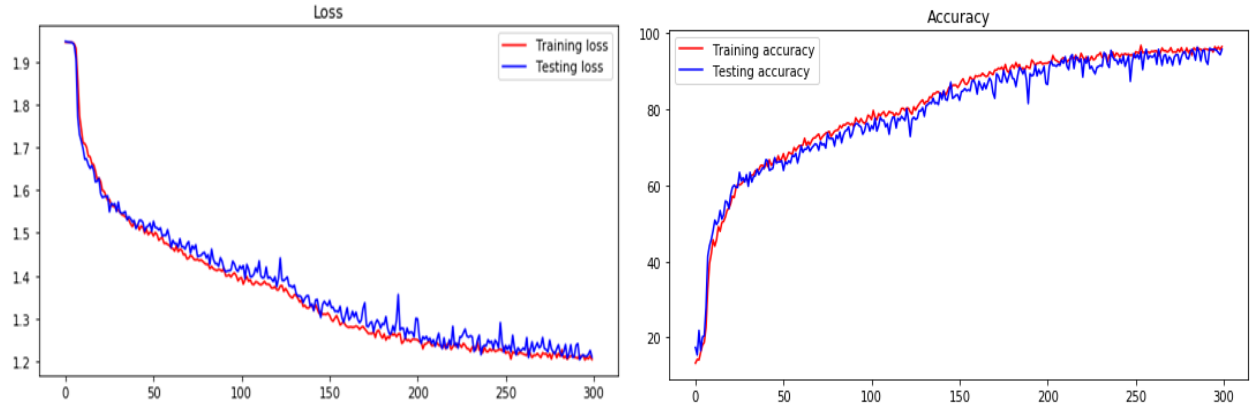


Figure 10: The variation trend of loss and accuracy of final version in the training process

```

RandomForestClassifier
When using RandomForest:
When k = 10
When use cross_val_score:
Accuracy: 0.98687(+/-)0.00003
When use cross_val_predict:
MSE: 0.09965
total time: 8.845

```

Figure 11: The training result using Random Forest

From the result in Figure 11 we can see that when use 10-fold cross validation, the accuracy is nearly 98% which is higher than our second net. However, when we used the high accuracy new net to do real time prediction, it performed badly so we finally didn't use the new net.

5 Source code

5.1 Dataset Preprocessing Using Openpose

Python file `preprocessing_using_openpose.py` is used to transfer the original video into a series of body key points. This file first transfer the video into frames and then input each frame into Openpose pretrained model to get the position of key points.

Note that if you want to run this file, you first need to have the pretrained model, which will be handed in to TA. As for our own dataset, the main part of Openpose is the same, so we won't list it here.

```
1 import cv2
2 import os
3 import random
4 import numpy as np
5 import scipy.io as io
6 import matplotlib.pyplot as plt
7
8 def get_path(path):
9     path_list = os.listdir(path)
10    path_list.sort()
11    full_path_list = []
12    for filename in path_list:
13        full_path_list.append(os.path.join(path, filename))
14    return full_path_list, len(full_path_list)
15
16 def get_frame_list(frame_number, get_number):
17     if frame_number < get_number:
18         flist = list(range(frame_number))
19         for k in range(get_number - frame_number):
20             flist.append(frame_number - 1)
21     elif frame_number < 1.3 * get_number:
22         flist = []
23         for k in range(get_number):
24             flist.insert(0, range(frame_number)[-1 - k])
25     else:
26         frame_list = list(range(round(frame_number * 0.18), frame_number))
27         flist = random.sample(frame_list, get_number)
28         flist.sort()
29
30     return flist
31
32 # cut the video into frames and implement openpose
33 def cut_video(path, path_folder_naming, arr_mat, arr_lab):
34     protoFile = r'C:\Users\Administrator\Desktop\AI_Project\openpose-master\models\pose\mpi\
35         pose_deploy_linevec.prototxt'
36     weightsFile = r'C:\Users\Administrator\Desktop\AI_Project\openpose-master\models\pose\
37         mpi\pose_iter_160000.caffemodel'
38     net = cv2.dnn.readNetFromCaffe(protoFile, weightsFile)
39     temp_list_x = []
40     temp_list_y = []
41
42     cap = cv2.VideoCapture(path)
43     category = eval(path[-5])
44     path_folder_naming[category - 1] += 1
45     naming_num = path_folder_naming[category - 1]
46     writing_folder = 'C:\\Users\\Administrator\\Desktop\\AI_Project\\dataset\\' + str(
47         category) + '\\\\' + str(naming_num) + '\\\\'
48     os.mkdir(writing_folder)
```

```

49     total_frame = int(cap.get(cv2.CAP_PROP_FRAMECOUNT))
    get_number = 8
    frame_list = get_frame_list(total_frame, get_number)
51
53     if cap.isOpened():
        print('video successfully opened')
55         success = True
57
    for i in range(get_number):
        cap.set(cv2.CAP_PROP_POS_FRAMES, frame_list[i])
59         success, frame = cap.read()
61
        print('Reading a new frame: ', success)
        if success:
63             cv2.imwrite(writing_folder + "frame" + "_%d.jpg" % i, frame, [int(cv2.
                IMWRITE_JPEG_QUALITY), 100])
            temp_list_x, temp_list_y = get_key_point_from_frame(frame, net, temp_list_x,
                temp_list_y)
65
        temp_list_x.extend(temp_list_y)
67         arr_mat.append(temp_list_x)
        arr_lab.append(category)
69
    cap.release()
71     return path_folder_naming, arr_mat, arr_lab
73
def get_key_point_from_frame(frame, net, temp_list_x, temp_list_y):
75
    # Specify number of points in the model
77     nPoints = 15
    im = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
79     inWidth = im.shape[1]
    inHeight = im.shape[0]
81
    # Convert image to blob
83     netInputSize = (368, 368)
    inpBlob = cv2.dnn.blobFromImage(im, 1.0 / 255, netInputSize, (0, 0, 0), swapRB=True,
        crop=False)
85     net.setInput(inpBlob)
87
    # Run Inference (forward pass)
    output = net.forward()
89     scaleX = float(inWidth) / output.shape[3]
    scaleY = float(inHeight) / output.shape[2]
91
    points = []
93     # Confidence threshold
    threshold = 0.1
95
    for i in range(nPoints):
97         probMap = output[0, i, :, :]
        minVal, prob, minLoc, point = cv2.minMaxLoc(probMap)
99
        # Scale the point to fit on the original image
101         x = scaleX * point[0]
        y = scaleY * point[1]
103
        if prob > threshold :
105             points.append((int(x), int(y)))
            temp_list_x.append(int(x))
107             temp_list_y.append(int(y))
        else :
109             points.append(None)

```

```

111         temp_list_x.append(None)
            temp_list_y.append(None)
113     return temp_list_x, temp_list_y
115 if __name__ == '__main__':
    folder_path = r'C:\Users\Administrator\Desktop\AI_Project\dataset\Florence_3d_actions'
117    video_path_list, video_num = get_path(folder_path)
    path_folder_naming = [0, 0, 0, 0, 0, 0, 0, 0, 0]
119    arr_mat = []
    arr_lab = []
121
    for num in range(video_num):
123        print('Now slice video', num + 1)
        path_folder_naming, arr_mat, arr_lab = cut_video(video_path_list[num],
            path_folder_naming, arr_mat, arr_lab)
125        print(path_folder_naming)
127
        if num == 100:
            # print(arr_mat)
129            print(np.array(arr_mat).shape)
            print(np.array(arr_lab).shape)
131
            np.save('sample_' + str(num) + '.npy', np.array(arr_mat))
133            np.save('label_' + str(num) + '.npy', np.array(arr_lab))
135
    arr_mat = np.array(arr_mat)
    arr_lab = np.array(arr_lab)
137    np.save('sample.npy', arr_mat)
    np.save('label.npy', arr_lab)

```

5.2 Training Process

Python file C2D_model.py is the code that we construct the 2D convolution training net.

The input shape of the net is (2,16,15) where 2 is the number of channels, 16 is the number of frames and 15 is the number of key points. The output is an array with 9 probs where each stands for the prediction probability of relative action.

```

import torch.nn as nn
2
class C2D(nn.Module):
4
    def __init__(self):
6        super(C2D, self).__init__()
8
        #input shape: (16, 15)
        self.conv1 = nn.Conv2d(2, 4, kernel_size=(3,3), padding=(1,1), stride=(1,1))#(16,
            15)
10        self.conv2 = nn.Conv2d(4, 8, kernel_size=(3,3), padding=(1,1), stride=(1,1))#(16,
            15)
        self.pool1 = nn.MaxPool2d(kernel_size=(2,2), stride=(2,2))#(8, 7)
12
        self.conv3 = nn.Conv2d(8, 12, kernel_size=(3,3), padding=(1,1), stride=(1,1))#(8, 7)
14        self.conv4 = nn.Conv2d(12,16, kernel_size=(3,3), padding=(1,1), stride=(1,1))#(8, 7)
        self.pool2 = nn.MaxPool2d(kernel_size=(2,2), stride=(2,2))#(4, 3)
16
        self.conv5 = nn.Conv2d(16,20, kernel_size=(3,3), padding=(1,1), stride=(1,1))#(4, 3)
18        self.conv6 = nn.Conv2d(20,24, kernel_size=(3,3), padding=(1,1), stride=(1,1))#(4, 3)
        self.pool3 = nn.MaxPool2d(kernel_size=(2,1), stride=(2,1))#(2, 3)
20
        self.fc1 = nn.Linear(144, 72)

```

```

22     self.fc2 = nn.Linear(72, 36)
23     self.fc3 = nn.Linear(36, 18)
24     self.fc4 = nn.Linear(18, 9)
25     #output labels: 9
26
27     self.dropout1 = nn.Dropout(p=0.1)
28     self.dropout2 = nn.Dropout(p=0.2)
29     self.dropout3 = nn.Dropout(p=0.3)
30     self.dropout4 = nn.Dropout(p=0.4)
31     self.dropout5 = nn.Dropout(p=0.5)
32
33     self.relu = nn.ReLU()
34     self.softmax = nn.Softmax()
35
36     def forward(self, x):
37
38         #convolution and pooling layers
39         h = self.relu(self.conv1(x))
40         h = self.relu(self.conv2(h))
41         h = self.pool1(h)
42
43         h = self.relu(self.conv3(h))
44         h = self.relu(self.conv4(h))
45         h = self.pool2(h)
46
47         h = self.relu(self.conv5(h))
48         h = self.relu(self.conv6(h))
49         h = self.pool3(h)
50
51         #flatten layer
52         h = h.view(-1, 144)
53
54         #full connection layers
55         h = self.relu(self.fc1(h))
56         h = self.dropout3(h)
57         h = self.relu(self.fc2(h))
58         h = self.dropout3(h)
59         h = self.relu(self.fc3(h))
60         h = self.dropout2(h)
61         logits = self.fc4(h)
62
63         probs = self.softmax(logits)
64
65         return probs

```

Python file `predict.py` is the code that we use to train the 2D convolution net and save the model we needed.

After loading and shuffling the key points dataset, we choose former 80% samples to be training set and later 20% samples to be training set. In the training process, we tried many combinations of parameters and finally got not bad results.

```

1  import torch
2  import torch.nn as nn
3
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import pandas as pd
7
8  from C2D_model import C2D
9
10 #calculate the accuracy
11 #y is actual label and z is prediction result(haven't changed to label)
12 def accuracy_cal(y, z):

```



```

13     batch_acc = 0
14     z_label = torch.argmax(z,1).numpy()
15     for m in range(len(y)):
16         if int(y[m]) == int(z_label[m]):
17             batch_acc += 1
18     return batch_acc
19
20 #testing set
21 def testing(net, testing_samples, testing_labels, loss_fn):
22     test_size = 1
23     test_iterate = int(len(testing_samples) / test_size)
24     test_acc = 0
25     test_loss = 0
26     for m in range(test_iterate):
27         y_pred = net(testing_samples[m*test_size:(m+1)*test_size])
28         y = testing_labels[m*test_size:(m+1)*test_size]
29         #calculate testing loss
30         loss = loss_fn(y_pred, y.type(torch.long))
31         test_loss += loss.item()
32         #calculate testing accuracy
33         test_acc += accuracy_cal(y,y_pred)
34
35     return test_loss, test_acc
36
37 #training process
38 def training_process(net, training_samples, training_labels, testing_samples, testing_labels,
39                     max_iterations, batch_size, learning_rate):
40
41     #set the parameters, batch size, max_iterations, test_size
42     batch_size = batch_size
43     max_iterations = max_iterations #200
44     iters_per_iterate = int(np.ceil(len(training_samples)/batch_size))
45
46     #define loss function
47     loss_fn = nn.CrossEntropyLoss(reduction = 'sum')
48
49     #define optimizer
50     optimizer = torch.optim.Adam(net.parameters(), lr = learning_rate) #0.001
51
52     #####
53     #begin training#
54     #####
55     #store the training and testing accuracy and loss information
56     training_loss, testing_loss = [], []
57     training_accuracy, testing_accuracy = [], []
58
59     for iterate in range(max_iterations):
60         iterate_loss = 0
61         iterate_acc = 0
62         for i in range(iters_per_iterate):
63             #if come to the end
64             if i == (iters_per_iterate -1) :
65                 x = (training_samples[i*batch_size: len(training_samples)])
66                 y = training_labels[i*batch_size: len(training_labels)]
67             else:
68                 x = (training_samples[i*batch_size: (i+1)*batch_size])
69                 y = training_labels[i*batch_size: (i+1)*batch_size]
70                 z = net(x)
71
72             #calculate the batch loss
73             loss_val = loss_fn(z, (y.type(torch.long)))
74             iterate_loss += loss_val.item()
75             #calculate the batch accuracy
76             iterate_acc += accuracy_cal(y,z)

```

```

77         #update the parameters
       optimizer.zero_grad()
79         loss_val.backward()
       optimizer.step()

81
       #update the acc and loss list
83       training_loss.append(iterate_loss/len(training_samples))
       training_accuracy.append(100 * iterate_acc / len(training_samples))
85
       #fit testing set
87       test_loss, test_acc = testing(net, testing_samples, testing_labels, loss_fn)
       testing_loss.append(test_loss / len(testing_samples))
89       testing_accuracy.append(100 * test_acc / len(testing_samples))
       print(iterate, training_accuracy[-1], testing_accuracy[-1])
91
       #plot the images
93       plot_image(training_loss, training_accuracy, testing_loss, testing_accuracy)
       #find the best model's number
95       print('The best model is in:', np.argmax(np.array(testing_accuracy)))
       return np.mean(np.array(training_accuracy[-1])), np.mean(np.array(testing_accuracy[-1]))
97
#plot the image of training and testing accuracy and loss in the training process
99 def plot_image(training_loss, training_accuracy, testing_loss, testing_accuracy):

101     figure, ax = plt.subplots(figsize = [9,10])
    x= range(len(training_loss))
103
    #plot the loss
105     plt.subplot(2,1,1)
    line1, =plt.plot(x, training_loss, 'r-', label = "Training loss")
107     plt.title('Loss')

109     plt.subplot(2,1,1)
    line2, =plt.plot(x, testing_loss, 'b-', label = 'Testing loss')
111     plt.legend(handles =[line1, line2], loc = 0)

113     #plot the accuracy
    plt.subplot(2,1,2)
115     line3, =plt.plot(x, training_accuracy, 'r-', label = 'Training accuracy')
    plt.title('Accuracy')
117
    plt.subplot(2,1,2)
119     line4, =plt.plot(x, testing_accuracy, 'b-', label = 'Testing accuracy')
    plt.legend(handles =[line3, line4], loc = 0)
121     #show the figure
    plt.show()
123
125
def training(net, features, labels, max_iterations, batch_size, learning_rate):
127
    #N: the number of samples, C: 2 channels (x,y), F: frames, H: height, actually 15 key
    points
129     N, C, F, H = features.shape

131     #reorder the trainig set
    np.random.seed(6)
133     np.random.shuffle(features)
    np.random.seed(6)
135     np.random.shuffle(labels)

137     #define the number of training set and testing set
    training_num = int(N * 0.8)
139     #get the training set and testing set
    training_samples = torch.Tensor(features[:training_num])

```

```

141 training_labels = torch.Tensor(labels[:training_num])
142 testing_samples = torch.Tensor(features[training_num:])
143 testing_labels = torch.Tensor(labels[training_num:])
144 #train the net
145 train_acc, test_acc = training_process(net, training_samples, training_labels,
146                                     testing_samples, testing_labels, max_iterations, batch_size, learning_rate)
147
148 #show the final accuracy
149 training_acc, testing_acc = [], []
150 training_acc.append(train_acc)
151 testing_acc.append(test_acc)
152 show = pd.DataFrame(columns = ('train_acc', 'test_acc'))
153 show['train_acc'] = training_acc
154 show['test_acc'] = testing_acc
155 print(show)
156
157 #save the model
158 torch.save(net, 'net.pkl')
159
160 #pre-treat the features and labels
161 def pre_treat(features, labels): #pre-treat the data
162
163     #delete some labels samples(which have many None data)
164     del_row = []
165     for i in range(len(labels)):
166         if int(labels[i]) in [4]:
167             del_row.append(i)
168     features = np.delete(features, del_row, 0)
169     labels = np.delete(labels, del_row, 0)
170
171     #change the remain labels to right order
172     for i in range(len(labels)):
173         if int(labels[i]) >= 1 and int(labels[i]) <= 3:
174             labels[i] = int(labels[i]) - 1
175         else:
176             labels[i] = int(labels[i]) - 2
177
178     return features, labels
179
180 #using value iteration to process None data
181 def filter_none_data(features):
182
183     #get the None data's position and set them to 0 initially
184     for i in range(len(features)):
185         non_pos = []
186         for j in range(len(features[i])):
187             if str(features[i][j]) == str(None):
188                 non_pos.append(j)
189                 features[i][j] = 0
190
191     #using 10 times value iteration to padding the missing values
192     for k in range(10): #value iteration times
193         for m in range(len(non_pos)):
194             if non_pos[m]%15 == 0: #begin position
195                 features[i][non_pos[m]] = features[i][non_pos[m]+1]
196             elif (non_pos[m]+1) % 15 == 0: #end position
197                 features[i][non_pos[m]] = features[i][non_pos[m]-1]
198             else: #medium position
199                 features[i][non_pos[m]] = (features[i][non_pos[m]+1] + features[i][
200                                     non_pos[m]-1]) / 2
201
202     return features
203
204 if __name__ == '__main__':

```

```

205     #load data
    features = np.load('sample_big.npy')
    labels = np.load('label_big.npy')
207
    #preprocess the features and labels
209     features, labels = pre_treat(features, labels)
    features = filter_none_data(features)
211
    #change the features to appropriate type
213     features = features.astype(np.int32)
215
    #N,C,F,H,W, actually no W
    features = features.reshape(-1,2,16,15)
217
    #initial the net
219     net = C2D()
221
    #assign the tuning parameters
    max_iterations = 300
223     batch_size = 4
    learning_rate = 0.0001
225
    #train the net parameters
227     training(net, features, labels, max_iterations, batch_size, learning_rate)
    print('max_iterations\t', 'batch_size\t', 'learning_rate')
229     print(max_iterations, '\t', batch_size, '\t', learning_rate)

```

5.3 Create Dataset Action-209

Python file `record.py` is the code we used to collect our dataset Action-209. Since we have color frames of Kinect which is just the real-time images it gets, we can just store them as the datasets. So we code it based on a previous demo. The only change is that we store images in a fixed frequency and write them to disks as the game runs.

During running the game, we can get the color frames which contains rgb-d information. We reshape them to image matrixes and split rgb and depth. Both of them are resized into (600, 480) and saved. The code is listed below.

```

from pykinect2 import PyKinectV2
2 from pykinect2.PyKinectV2 import *
from pykinect2 import PyKinectRuntime
4
6 import ctypes
import _ctypes
8 import pygame
import sys
10 import math
import numpy as np
12 import cv2
import time
14 import os
16
'''
18 This py code is used to collect our dataset containing rgb images and depth data
it is based on a demo of drawing box (not the final version)
the function is done at # 310~340, So I only wrote notation at that part
20 the full notation is at done in the final virsion demo
'''
22

```

```

24 if sys.hexversion >= 0x03000000:
    import _thread as thread
else:
26     import thread

28 # path to store the data
PATH = "E:/AI/final/PyKinect2/examples/video/10/3/"

30 # colors for drawing different bodies
32 SKELETON_COLORS = [pygame.color.THECOLORS["red"],
    pygame.color.THECOLORS["blue"],
34     pygame.color.THECOLORS["green"],
    pygame.color.THECOLORS["orange"],
36     pygame.color.THECOLORS["purple"],
    pygame.color.THECOLORS["yellow"],
38     pygame.color.THECOLORS["violet"]]

40
class BodyGameRuntime(object):
42     def __init__(self):
        pygame.init()

44
        # Used to manage how fast the screen updates
46         self._clock = pygame.time.Clock()

48         # Set the width and height of the screen [width, height]
        self._infoObject = pygame.display.Info()
50         self._screen = pygame.display.set_mode((self._infoObject.current_w >> 1, self.
            _infoObject.current_h >> 1),
            pygame.HWSURFACE|pygame.DOUBLEBUF|pygame.
52                 RESIZABLE, 32)

        pygame.display.set_caption("Action Recognition with Kinect")

54
        # Loop until the user clicks the close button.
56         self._done = False

58         # Used to manage how fast the screen updates
        self._clock = pygame.time.Clock()

60
        # Kinect runtime object, we want only color and body frames
62         self._kinect = PyKinectRuntime.PyKinectRuntime(PyKinectV2.FrameSourceTypes_Color |
            PyKinectV2.FrameSourceTypes_Body)

64         # back buffer surface for getting Kinect color frames, 32bit color, width and height
            equal to the Kinect color frame size
        self._frame_surface = pygame.Surface((self._kinect.color_frame_desc.Width, self.
66             _kinect.color_frame_desc.Height), 0, 32)

        # here we will store skeleton data
68         self._bodies = None

70         # pose of user
        self._pose = "stand"

72
        #to get data
74         self._is_get = False

76         #begin record
        self._frame_num = 0
78         if os.listdir(PATH) == []:
            self._group_num = 0
80         else:
            n = 0
82             l = os.listdir(PATH)

```

```

84         for i in l:
85             m = int(i)
86             if m > n:
87                 n = m
88             self._group_num = n + 1
89
90         self.depth = np.zeros((1,600,480))
91         # data
92         self._x = []
93         self._y = []
94         self.labels = ["wave", "applause", "sit", "stand", "watch wrist"]
95
96         # interval
97         self._inter_range = 100
98         self._interval = 0
99
100         print("model loaded")
101
102     def draw_body_bone(self, joints, jointPoints, color, joint0, joint1):
103         joint0State = joints[joint0].TrackingState;
104         joint1State = joints[joint1].TrackingState;
105
106         # both joints are not tracked
107         if (joint0State == PyKinectV2.TrackingState_NotTracked) or (joint1State ==
108             PyKinectV2.TrackingState_NotTracked):
109             return
110
111         # both joints are not *really* tracked
112         if (joint0State == PyKinectV2.TrackingState_Inferred) and (joint1State == PyKinectV2
113             .TrackingState_Inferred):
114             return
115
116         # ok, at least one is good
117         start = (jointPoints[joint0].x, jointPoints[joint0].y)
118         end = (jointPoints[joint1].x, jointPoints[joint1].y)
119
120         try:
121             pygame.draw.line(self._frame_surface, color, start, end, 8)
122         except: # need to catch it due to possible invalid positions (with inf)
123             pass
124
125     def draw_body(self, joints, jointPoints, color):
126         # Torso
127         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_Head,
128             PyKinectV2.JointType_Neck);
129         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_Neck,
130             PyKinectV2.JointType_SpineShoulder);
131         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineShoulder,
132             PyKinectV2.JointType_SpineMid);
133         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineMid,
134             PyKinectV2.JointType_SpineBase);
135         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineShoulder,
136             PyKinectV2.JointType_ShoulderRight);
137         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineShoulder,
138             PyKinectV2.JointType_ShoulderLeft);
139         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineBase,
140             PyKinectV2.JointType_HipRight);
141         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineBase,
142             PyKinectV2.JointType_HipLeft);
143
144         # Right Arm
145         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ShoulderRight,
146             PyKinectV2.JointType_ElbowRight);
147         self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ElbowRight,
148             PyKinectV2.JointType_WristRight);

```

```

136     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristRight,
        PyKinectV2.JointType_HandRight);
137     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HandRight,
        PyKinectV2.JointType_HandTipRight);
138     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristRight,
        PyKinectV2.JointType_ThumbRight);

140     # Left Arm
        self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ShoulderLeft,
            PyKinectV2.JointType_ElbowLeft);
142     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ElbowLeft,
        PyKinectV2.JointType_WristLeft);
        self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristLeft,
            PyKinectV2.JointType_HandLeft);
144     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HandLeft,
        PyKinectV2.JointType_HandTipLeft);
        self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristLeft,
            PyKinectV2.JointType_ThumbLeft);

146     # Right Leg
148     #self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HipRight,
        PyKinectV2.JointType_KneeRight);
        self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_KneeRight,
            PyKinectV2.JointType_AnkleRight);
150     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_AnkleRight,
        PyKinectV2.JointType_FootRight);

152     # Left Leg
        self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HipLeft,
            PyKinectV2.JointType_KneeLeft);
154     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_KneeLeft,
        PyKinectV2.JointType_AnkleLeft);
        self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_AnkleLeft,
            PyKinectV2.JointType_FootLeft);

156

158     def draw_box(self, joints, jointPoints, color):
159         X = []
160         Y = []
        self._interval += self._clock.get_time()
162         if self._interval > self._inter_range:
            self._interval = 0
            self._is_get = True
164         # [3] [20] [8] [9] [10] [4] [5] [6] [16] [17] [18] [12] [13] [14] [1]
166         points = [PyKinectV2.JointType_Head, PyKinectV2.JointType_SpineShoulder, PyKinectV2.
            JointType_ShoulderRight,
                PyKinectV2.JointType_ElbowRight, PyKinectV2.JointType_WristRight, PyKinectV2.
                    JointType_ShoulderLeft,
168                PyKinectV2.JointType_ElbowLeft, PyKinectV2.JointType_WristLeft, PyKinectV2.
                    JointType_HipRight,
                PyKinectV2.JointType_KneeRight, PyKinectV2.JointType_AnkleRight, PyKinectV2.
                    JointType_HipLeft,
170                PyKinectV2.JointType_KneeLeft, PyKinectV2.JointType_AnkleLeft, PyKinectV2.
                    JointType_SpineMid,
                PyKinectV2.JointType_Neck, PyKinectV2.JointType_SpineBase, PyKinectV2.
                    JointType_HandRight,
172                PyKinectV2.JointType_HandTipRight, PyKinectV2.JointType_ThumbRight,
                PyKinectV2.JointType_HandLeft,
                PyKinectV2.JointType_HandTipLeft, PyKinectV2.JointType_ThumbLeft, PyKinectV2.
                    JointType_FootRight,
174                PyKinectV2.JointType_FootLeft]
        left_most = float('inf')
176        right_most = 0
        up_most = float('inf')
178        down_most = 0

```

```

180     for No in range(25):
181         pt = points[No]
182         jointState = joints[pt].TrackingState;
183         if jointState == PyKinectV2.TrackingState_NotTracked or jointState == PyKinectV2
184             .TrackingState_Inferred:
185             if self._is_get == True and No < 15:
186                 X.append(None)
187                 Y.append(None)
188             continue
189             x = jointPoints[pt].x
190             y = jointPoints[pt].y
191             if x < left_most:
192                 left_most = x
193             if x > right_most:
194                 right_most = x
195             if y < up_most:
196                 up_most = y
197             if y > down_most:
198                 down_most = y
199             if self._is_get == True and No < 15:
200                 # print(No)
201                 # print(pt)
202                 # print(x, y)
203                 if x < 0 or x > 1920:
204                     X.append(None)
205                 else:
206                     X.append(int(x * 600 / 1920))
207                     # self._is_get = False
208                     # self._interval = self._inter_range
209                 if y < 0 or y > 1080:
210                     Y.append(None)
211                 else:
212                     Y.append(int(y * 480 / 1080))
213                     # self._is_get = False
214                     # self._interval = self._inter_range
215
216             if self._is_get:
217                 if len(X)==15 and len(Y)==15:
218                     self._x += X
219                     self._y += Y
220                 else:
221                     self._interval = self._inter_range
222                     self._is_get = False
223                     # print(len(self._x), len(self._y))
224                     # print(self._x)
225                     # print(self._y)
226             height = down_most - up_most
227             width = right_most - left_most
228
229             corner1 = (left_most - 0.7*math.sqrt(abs(width)), up_most - 0.15*height)
230             corner2 = (right_most + 0.7*math.sqrt(abs(width)), up_most - 0.15*height)
231             corner3 = (right_most + 0.7*math.sqrt(abs(width)), down_most + 0.15*height)
232             corner4 = (left_most - 0.7*math.sqrt(abs(width)), down_most + 0.15*height)
233             try:
234                 pygame.draw.line(self._frame_surface, color, corner1, corner2, 8)
235                 pygame.draw.line(self._frame_surface, color, corner2, corner3, 8)
236                 pygame.draw.line(self._frame_surface, color, corner3, corner4, 8)
237                 pygame.draw.line(self._frame_surface, color, corner4, corner1, 8)
238
239             except: # need to catch it due to possible invalid positions (with inf)
240                 pass
241
242     def update_pose(self):
243         if len(self._x) < 120 or len(self._y) < 120:

```



```

244         return
245         if len(self._x) > 120 or len(self._y) > 120:
246             print("range error!")
247             self._x = []
248             self._y = []
249         return
250         data = np.array(self._x + self._y)
251         data = np.reshape(data, (1, 240))
252         self._x = []
253         self._y = []
254
255         data = data.reshape(-1, 2, 8, 15)
256         data = torch.Tensor(data)
257
258         predict = np.argmax(self.net(data).data.numpy())
259         self._pose = self.labels[predict]
260
261         self._interval = 0
262
263     def draw_color_frame(self, frame, target_surface):
264         target_surface.lock()
265         address = self._kinect.surface_as_array(target_surface.get_buffer())
266         ctypes.memmove(address, frame.ctypes.data, frame.size)
267         del address
268         target_surface.unlock()
269
270     def run(self):
271         # ----- Main Program Loop -----
272         while not self._done:
273             # --- Main event loop
274             for event in pygame.event.get(): # User did something
275                 if event.type == pygame.QUIT: # If user clicked close
276                     self._done = True # Flag that we are done so we exit this loop
277
278                 elif event.type == pygame.VIDEORESIZE: # window resized
279                     self._screen = pygame.display.set_mode(event.dict['size'],
280                                                             pygame.HWSURFACE|pygame.DOUBLEBUF|pygame.
281                                                             RESIZABLE, 32)
282
283             # --- Game logic should go here
284
285             # --- Getting frames and drawing
286
287             # --- Cool! We have a body frame, so can get skeletons
288             if self._kinect.has_new_body_frame():
289                 self._bodies = self._kinect.get_last_body_frame()
290
291             # --- draw skeletons to _frame_surface
292             if self._bodies is not None:
293                 for i in range(0, self._kinect.max_body_count):
294                     body = self._bodies.bodies[i]
295                     if not body.is_tracked:
296                         continue
297
298                     joints = body.joints
299                     # convert joint coordinates to color space
300                     joint_points = self._kinect.body_joints_to_color_space(joints)
301                     #self.draw_body(joints, joint_points, SKELETON_COLORS[i])
302
303             # --- Woohoo! We've got a color frame! Let's fill out back buffer surface with
304             # frame's data
305             if self._kinect.has_new_color_frame():
306                 self._interval += self._clock.get_time()

```

```

306         if self._interval > self._inter_range:
307             self._interval = 0
308             self._is_get = True
309
310
311         '''
312         here is the main code to collect the dataset
313         '''
314         frame = self._kinect.get_last_color_frame()
315         # reshape the frame into image shape
316         image = frame.reshape((self._kinect.color_frame_desc.Height, self._kinect.
317                                color_frame_desc.Width, 4), order='C')
318         # get the rgb part and shrink them
319         rgb_image = image[:, :, :3].copy()
320         rgb_image = cv2.resize(rgb_image, (600, 480))
321         # get the depth part and shrink it
322         d_image = image[:, :, 3].copy()
323         d_image = cv2.resize(d_image, (600, 480))
324         d_image = np.reshape(d_image, (1,600,480))
325
326         if self._is_get:
327             if self._frame_num == 0:
328                 os.mkdir(PATH+str(self._group_num)+'/')
329                 path = PATH + str(self._group_num) + '/'
330                 # save the rgb image
331                 cv2.imwrite(path+str(self._frame_num)+".jpg", rgb_image)
332                 # concatenate 16 depth data together
333                 self.depth = np.concatenate((self.depth, d_image), axis=0)
334                 self._frame_num += 1
335                 if self._frame_num == 16:
336                     self.depth = np.delete(self.depth, 0, axis=0)
337                     # save all the 16 frames' depth data
338                     np.save(path+"depth.npy", self.depth)
339                     self.depth = np.zeros((1,600,480))
340                     self._frame_num = 0
341                     self._group_num += 1
342
343
344
345         # cv2.imshow('test', image)
346         self.draw_color_frame(frame, self._frame_surface)
347         frame = None
348
349
350
351         # —— copy back buffer surface pixels to the screen, resize it if needed and
352         # —— keep aspect ratio
353         # —— (screen size may be different from Kinect's color frame size)
354
355         h_to_w = float(self._frame_surface.get_height()) / self._frame_surface.get_width()
356         target_height = int(h_to_w * self._screen.get_width())
357         surface_to_draw = pygame.transform.scale(self._frame_surface, (self._screen.
358                                                    get_width(), target_height))
359
360         ## Display some text
361         # font = pygame.font.Font(None, 60)
362         # text = font.render(self._pose, 1, (10, 10, 10))
363         # textpos = text.get_rect()
364         # textpos.centerx = surface_to_draw.get_rect().centerx
365         # surface_to_draw.blit(text, textpos)
366
367         self._screen.blit(surface_to_draw, (0,0))
368         pygame.display.update()

```

```

368         # — Go ahead and update the screen with what we've drawn.
        pygame.display.flip()
370
372         # — Limit to 60 frames per second
        self._clock.tick(60)
374
376         # Close our Kinect sensor, close the window and quit.
        self._kinect.close()
        pygame.quit()
378
380 __main__ = "Kinect v2 Body Game"
game = BodyGameRuntime();
game.run();

```

5.4 Real-time Recognition

Python file `PyKinectBodyGame.py` is the demo of real-time action recognition system. It is a game based on `Pygame`. After running it, a game window will appear. Stand in front of it, you will see a box encircling you, which means you are sensed. After that, do some actions in the list ['wave', 'drink', 'call', 'applause', 'stand', 'sit', 'stand still'], and the label it recognized will be shown on the screen.

```

from pykinect2 import PyKinectV2
2 from pykinect2.PyKinectV2 import *
from pykinect2 import PyKinectRuntime
4
from C2D_model import C2D
6 import torch
import torch.nn as nn
8 from torch.autograd import Variable

10 import ctypes
import _ctypes
12 import pygame
import sys
14 import math
import numpy as np
16 import cv2
import time
18 import torch

20 if sys.hexversion >= 0x03000000:
    import _thread as thread
22 else:
    import thread
24

# colors for drawing different bodies
26 SKELETON.COLORS = [pygame.color.THECOLORS["red"],
    pygame.color.THECOLORS["blue"],
28     pygame.color.THECOLORS["green"],
    pygame.color.THECOLORS["orange"],
30     pygame.color.THECOLORS["purple"],
    pygame.color.THECOLORS["yellow"],
32     pygame.color.THECOLORS["violet"]]

34
class BodyGameRuntime(object):
36     def __init__(self):

```

```

38     # initialize the game
    pygame.init()

40     # Used to manage how fast the screen updates
    self._clock = pygame.time.Clock()

42     # Set the width and height of the screen [width, height]
44     self._infoObject = pygame.display.Info()
    self._screen = pygame.display.set_mode((self._infoObject.current_w >> 1, self.
        _infoObject.current_h >> 1),
46                                         pygame.HWSURFACE|pygame.DOUBLEBUF|pygame.
                                            RESIZABLE, 32)

48     pygame.display.set_caption("Action Recognition with Kinect")

50     # Loop until the user clicks the close button.
    self._done = False

52     # Used to manage how fast the screen updates
54     self._clock = pygame.time.Clock()

56     # Kinect runtime object, we want only color and body frames
    self._kinect = PyKinectRuntime.PyKinectRuntime(PyKinectV2.FrameSourceTypes_Color |
60         PyKinectV2.FrameSourceTypes_Body)

58     # back buffer surface for getting Kinect color frames, 32bit color, width and height
    # equal to the Kinect color frame size
60     self._frame_surface = pygame.Surface((self._kinect.color_frame_desc.Width, self.
        _kinect.color_frame_desc.Height), 0, 32)

62     # here we will store skeleton data
    self._bodies = None

64     # pose of user
66     self._pose = "stand"

68     #to get data
    self._is_get = False

70     # data
72     self._x = [] # x coordiates of joints points
    self._y = [] # y coordiates of joints points
74     self.labels = ['wave', 'drink', "call", "appaluse", "stand", "sit", "stand still"]

76     # interval
    self._inter_range = 100
78     self._interval = 0

80     # build the net by loading pytorch model
    self.net = torch.load('net2.pkl')
82     print("model loaded")

84     # show a bone on the screen
    def draw_body_bone(self, joints, jointPoints, color, joint0, joint1):
86         joint0State = joints[joint0].TrackingState;
        joint1State = joints[joint1].TrackingState;

88         # both joints are not tracked
90         if (joint0State == PyKinectV2.TrackingState_NotTracked) or (joint1State ==
            PyKinectV2.TrackingState_NotTracked):
            return

92         # both joints are not *really* tracked
94         if (joint0State == PyKinectV2.TrackingState_Inferred) and (joint1State == PyKinectV2
            .TrackingState_Inferred):

```

```

96         return
97
98     # get the endpoints
99     start = (jointPoints[joint0].x, jointPoints[joint0].y)
100     end = (jointPoints[joint1].x, jointPoints[joint1].y)
101
102     try:
103         pygame.draw.line(self._frame_surface, color, start, end, 8)
104     except: # need to catch it due to possible invalid positions (with inf)
105         pass
106
107 # draw the total body on the screen
108 def draw_body(self, joints, jointPoints, color):
109     # Torso
110     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_Head,
111                         PyKinectV2.JointType_Neck);
112     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_Neck,
113                         PyKinectV2.JointType_SpineShoulder);
114     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineShoulder,
115                         PyKinectV2.JointType_SpineMid);
116     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineMid,
117                         PyKinectV2.JointType_SpineBase);
118     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineShoulder,
119                         PyKinectV2.JointType_ShoulderRight);
120     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineShoulder,
121                         PyKinectV2.JointType_ShoulderLeft);
122     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineBase,
123                         PyKinectV2.JointType_HipRight);
124     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_SpineBase,
125                         PyKinectV2.JointType_HipLeft);
126
127     # Right Arm
128     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ShoulderRight,
129                         PyKinectV2.JointType_ElbowRight);
130     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ElbowRight,
131                         PyKinectV2.JointType_WristRight);
132     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristRight,
133                         PyKinectV2.JointType_HandRight);
134     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HandRight,
135                         PyKinectV2.JointType_HandTipRight);
136     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristRight,
137                         PyKinectV2.JointType_ThumbRight);
138
139     # Left Arm
140     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ShoulderLeft,
141                         PyKinectV2.JointType_ElbowLeft);
142     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_ElbowLeft,
143                         PyKinectV2.JointType_WristLeft);
144     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristLeft,
145                         PyKinectV2.JointType_HandLeft);
146     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HandLeft,
147                         PyKinectV2.JointType_HandTipLeft);
148     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_WristLeft,
149                         PyKinectV2.JointType_ThumbLeft);
150
151     # Right Leg
152     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HipRight,
153                         PyKinectV2.JointType_KneeRight);
154     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_KneeRight,
155                         PyKinectV2.JointType_AnkleRight);
156     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_AnkleRight,
157                         PyKinectV2.JointType_FootRight);
158
159     # Left Leg
160     self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_HipLeft,

```

```

        PyKinectV2.JointType_KneeLeft);
self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_KneeLeft,
PyKinectV2.JointType_AnkleLeft);
140 self.draw_body_bone(joints, jointPoints, color, PyKinectV2.JointType_AnkleLeft,
PyKinectV2.JointType_FootLeft);

142
# draw the box on the screen
144 def draw_box(self, joints, jointPoints, color):
'''
146     function1: draw the box of body, according to the boundary of coordinates
    function2: store the data if it is time to sample
148 '''
    X = []
    Y = []
    # calculate the interval from last sampling
    152 self._interval += self._clock.get_time()
    # if it is time to sample
    154 if self._interval > self._inter_range:
        self._interval = 0
        156 self._is_get = True
    # [3] [20] [8] [9] [10] [4] [5] [6] [16] [17] [18] [12] [13] [14] [1]
    158 # the first 15 joints are the data we need to input to the model
    # they have the same order with training datasets
    160 points = [PyKinectV2.JointType_Head, PyKinectV2.JointType_SpineShoulder, PyKinectV2.
        JointType_ShoulderRight,
        PyKinectV2.JointType_ElbowRight, PyKinectV2.JointType_WristRight, PyKinectV2.
        JointType_ShoulderLeft,
    162 PyKinectV2.JointType_ElbowLeft, PyKinectV2.JointType_WristLeft, PyKinectV2.
        JointType_HipRight,
        PyKinectV2.JointType_KneeRight, PyKinectV2.JointType_AnkleRight, PyKinectV2.
        JointType_HipLeft,
    164 PyKinectV2.JointType_KneeLeft, PyKinectV2.JointType_AnkleLeft, PyKinectV2.
        JointType_SpineMid,
        PyKinectV2.JointType_Neck, PyKinectV2.JointType_SpineBase, PyKinectV2.
        JointType_HandRight,
    166 PyKinectV2.JointType_HandTipRight, PyKinectV2.JointType_ThumbRight,
        PyKinectV2.JointType_HandLeft,
        PyKinectV2.JointType_HandTipLeft, PyKinectV2.JointType_ThumbLeft, PyKinectV2.
        JointType_FootRight,
    168 PyKinectV2.JointType_FootLeft]

    # initialize the boundary
    170 left_most = float('inf')
    right_most = 0
    172 up_most = float('inf')
    down_most = 0
    # traverse all the joints
    174 for No in range(25):
        pt = points[No]
        jointState = joints[pt].TrackingState
        178 # if the joint is not tracked(missing) and we need to get it
        # store the coordinates as None and continue
        180 if jointState == PyKinectV2.TrackingState_NotTracked or jointState == PyKinectV2.
            TrackingState_Inferred:
            if self._is_get == True and No < 15:
                182 X.append(None)
                Y.append(None)
            184 continue
        # update the boundary
        186 x = jointPoints[pt].x
        y = jointPoints[pt].y
        188 if x < left_most:
            left_most = x
        190 if x > right_most:
            right_most = x

```

```

192         if y < up_most:
193             up_most = y
194         if y > down_most:
195             down_most = y
196         # if it is time to sample
197         if self._is_get == True and No < 15:
198             # if out of range, append None
199             if x < 0 or x > 1920:
200                 X.append(None)
201             else:
202                 # change the size of coordinates and store
203                 X.append(600 - int(x * 600 / 1920))
204             if y < 0 or y > 1080:
205                 Y.append(None)
206             else:
207                 Y.append(int(y * 480 / 1080))
208
209         # append the data of this frame to the total list
210         if self._is_get:
211             if len(X)==15 and len(Y)==15:
212                 self._x += X
213                 self._y += Y
214             # if data is not complete, sample next frame
215             # but I think it is not gonna happen
216             else:
217                 self._interval = self._inter_range
218         height = down_most - up_most
219         width = right_most - left_most
220
221         # since the joints is not the strict boundary
222         # change the coordinates
223         corner1 = (left_most - 0.7*math.sqrt(abs(width)), up_most - 0.15*height)
224         corner2 = (right_most + 0.7*math.sqrt(abs(width)), up_most - 0.15*height)
225         corner3 = (right_most + 0.7*math.sqrt(abs(width)), down_most + 0.15*height)
226         corner4 = (left_most - 0.7*math.sqrt(abs(width)), down_most + 0.15*height)
227         # draw the box
228         try:
229             pygame.draw.line(self._frame_surface, color, corner1, corner2, 8)
230             pygame.draw.line(self._frame_surface, color, corner2, corner3, 8)
231             pygame.draw.line(self._frame_surface, color, corner3, corner4, 8)
232             pygame.draw.line(self._frame_surface, color, corner4, corner1, 8)
233
234         except: # need to catch it due to possible invalid positions (with inf)
235             pass
236
237         # filter the None data
238         '''
239         if we get None in draw_box()
240         set them with approximate values
241         '''
242         def filter_none_data(self, features):
243             for i in range(len(features)):
244                 # set None to 0
245                 non_pos = []
246                 for j in range(len(features[i])):
247                     if str(features[i][j]) == str(None):
248                         non_pos.append(j)
249                         features[i][j] = 0
250                 # if it is 0, set it with average of neighbors
251                 for k in range(10): #value iteration times
252                     for m in range(len(non_pos)):
253                         if non_pos[m]%15 == 0: #begin position
254                             features[i][non_pos[m]] = features[i][non_pos[m]+1]
255                         elif (non_pos[m]+1) % 15 == 0: #end position
256                             features[i][non_pos[m]] = features[i][non_pos[m]-1]

```

```

258         else: #medium position
                features[i][non_pos[m]] = (features[i][non_pos[m]+1] + features[i][
                    non_pos[m]-1]) / 2
        return features

260

262 def update_pose(self):
    '''
264     update the pose by feed the data into the model
    show the label in the screen
266     '''
    # if data is not enough, skip
268     if len(self._x) < 240 or len(self._y) < 240:
        return
    # if data out of range, raise error
270     if len(self._x) > 240 or len(self._y) > 240:
        raise ValueError("out of range")
    Time = time.time()
274     data = np.array(self._x + self._y)
    data = np.reshape(data, (1, 480))
276     data = self.filter_none_data(data)
    self._x = []
278     self._y = []

280     data = data.astype(np.int32)
    # reshape the data to fit the model
282     data = data.reshape(-1,2,16,15)
    # turn into torch tensor
284     data = torch.Tensor(data)

286     predict = np.argmax(self.net(data).data.numpy())
    # update the label
288     self._pose = self.labels[predict]
    # restart the interval
290     self._interval = 0
    # show the time of predict
292     print(time.time()-Time)

294

    # draw the frame on the screen
296 def draw_color_frame(self, frame, target_surface):
    target_surface.lock()
298     address = self._kinect.surface_as_array(target_surface.get_buffer())
    ctypes.memmove(address, frame.ctypes.data, frame.size)
300     del address
    target_surface.unlock()

302

    # run pygame
304 def run(self):
    # ----- Main Program Loop -----
306     while not self._done:
        # Main event loop
308         for event in pygame.event.get(): # User did something
            if event.type == pygame.QUIT: # If user clicked close
310                 self._done = True # Flag that we are done so we exit this loop

312         elif event.type == pygame.VIDEORESIZE: # window resized
            self._screen = pygame.display.set_mode(event.dict['size'],
314                                                     pygame.HWSURFACE|pygame.DOUBLEBUF|pygame.
                                                         RESIZABLE, 32)

316

        # Getting frames and drawing
318         if self._kinect.has_new_color_frame():
            frame = self._kinect.get_last_color_frame()

```



```

320         # it can be reshape into image matrix(RGB-D)
322         # image = frame.reshape((self._kinect.color_frame_desc.Height, self._kinect.
            color_frame_desc.width, 4),order='C')
324         # image = image[:, :, :3]
            # cv2.imshow('test', image)

326         self.draw_color_frame(frame, self._frame_surface)
            frame = None

328
330         # We have a body frame, so can get skeletons
            if self._kinect.has_new_body_frame():
332                 self._bodies = self._kinect.get_last_body_frame()

334         # draw skeletons to _frame_surface
            if self._bodies is not None:
336                 # traverse the bodies
                    for i in range(0, self._kinect.max_body_count):
338                         body = self._bodies.bodies[i]
                            if not body.is_tracked:
                                continue

340
342                         joints = body.joints
                            # convert joint coordinates to color space
                                joint_points = self._kinect.body_joints_to_color_space(joints)
344                                # draw the box and store the data
                                    self.draw_box(joints, joint_points, SKELETON_COLORS[i])

346
348         # copy back buffer surface pixels to the screen, resize it if needed and keep
            aspect_ratio
350         # (screen size may be different from Kinect's color frame size)
            self.update_pose()
352         h_to_w = float(self._frame_surface.get_height()) / self._frame_surface.get_width
            ()
            target_height = int(h_to_w * self._screen.get_width())
            surface_to_draw = pygame.transform.scale(self._frame_surface, (self._screen.
                get_width(), target_height))

354         # Display the action label
            font = pygame.font.Font(None, 60)
356         text = font.render(self._pose, 1, (255, 10, 10))
            textpos = text.get_rect()
358         textpos.centerx = surface_to_draw.get_rect().centerx
            surface_to_draw.blit(text, textpos)
360         self._screen.blit(surface_to_draw, (0,0))
            # update the frame
362         pygame.display.update()

364         # Go ahead and update the screen with what we've drawn.
            pygame.display.flip()

366
368         # Limit the fps
            self._clock.tick(20)

370         # Close our Kinect sensor, close the window and quit.
            self._kinect.close()
372         pygame.quit()

374
__main__ = "Kinect v2 Body Game"
376 game = BodyGameRuntime();
game.run();

```

6 Discussion and Conclusions

6.1 Discussion

6.1.1 Missing Data Processing

Openpose is based on RGB graph. In some condition it's not as powerful as we thought. For instance, if one's arm is hidden behind his/her body, Openpose cannot recognize, or very inaccurate. We thought of two solutions for this problem. First is to infer the missing point with its neighbors. Second is to implement depth graph. For the first solution, we just use 10 times value iteration to fill these values by neighboring ones. However, it's not so sensible because the relationship of these values is not evidently known by us. When we want to add boxing and walking to our actions, we found the net performs badly due to the roughly processing of None values. The second solution will be introduced later.

In the future, we may try to use HMM to help us with padding these values.

6.1.2 The Use of Depth Information

The Kinect can provide us with RGB-D images and skeleton structure. First, we use RGB images to train our 3D convolution model. Due to the large number of dimensions and many irrelevant and useless information, our first trying only attended nearly 25% accuracy. Then, we use the key points of the skeleton structure in our net mentioned above and obtained not bad results.

However, the remaining depth images we didn't use is also a key breakthrough point. We saved the depth images when we created our own dataset. At that time, because of running out of memory when we load the depth data, we didn't use it. Then we reduced the dimension of the depth images and changed the structure of our model to use both key points and depth information. What finally frustrated us was that when we were ready to load the small size depth data, we found that the depth data were all 255! However, we don't have enough time to capture the actions again. What a shame...

6.1.3 Our New Dataset

As for our new dataset. There are three aspects we want to improve.

First, we want to add high-quality depth images, shown in Figure 12. Although we collected depth images in the first vision, we found that the matrix contains 255 mostly, which means it is not useful while training. We believe it is because we stood too close to the wall at that time so that the Kinect could not distinguish the depth precisely. We plan to record the depth images one more time during which we will stand far from the background. But here comes another question, with different distance from the background we will get totally different data range in depth images. What if we have new testing situations with various distances from the background? Can our model learn the generalization ability training with monotonous background dataset? That's the thing we have to figure out if we use depth images to train our model.



Figure 12: We intend to implement depth graphs.

Second, we plan to expand our action pool with some interactive actions, for example, fighting and embracing (Figure 13). They are also very useful in real life application. It is also very challenging. How does our model know that it is time to recognize interactive actions rather than single actions of 2 people? What's more, in interactive actions, parts of players' bodies will sure to overlap. It is hard to get joints coordinates with such an incomplete information and it is also very hard to predict actions with them.

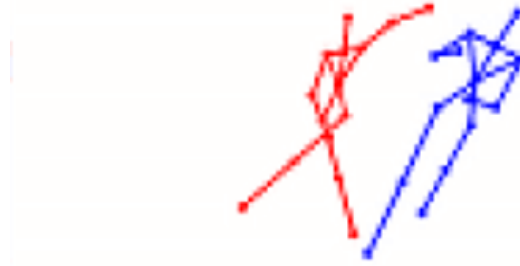


Figure 13: We intend to add some interactive actions.

Finally, we plan to expand the dataset with data containing more people. Limited by time and resources, we only collected data of the group members. Of course, it is not big enough because different people have different bodily form and action habits. Those are import part of the model generalization ability.

6.1.4 Real-time Video Sampling

We found that it still makes some mistakes during interval of actions. It is because the data during the interval of different actions is not related to any action. Of course, we did not train our model with this kind of data. It is hard to predict actions like that because they are not defined actions at all. We plan to handle it like this. We only show the action label when it is confirmed by our model. We can set the threshold value as 50%, only when the possibility is more than it can we confirm the prediction.

6.2 Conclusion

In this project, we presented a convolution network to perform human action recognition. We used OpenPose to extract joints coordinates which can show key information of a person's posture. We concatenated series of coordinates and input them to get prediction label. This CNN is first trained on FLORENCE 3D ACTIONS DATASET, which is very small with 8 frames formed as a sample. To improve the generalization ability of the CNN model we collected our own dataset called Action-209. We designed it according to the shortcomings of the previous one. Then we adjusted the structure CNN to fit the sample size of 16 frames and get testing accuracy at nearly 95%. The model is also applied on a real-time action recognition system. It can do multi-targets recognition at the same time and perform with considerable accuracy. That means our model does have practical application value.

References

- [1] Diogo C. Luvizon¹, David Picard, Hedi Tabia. “2D/3D Pose Estimation and Action Recognition using Multitask Deep Learning,” in CVPR, 2018
- [2] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. “Learning spatiotemporal features with 3d convolutional networks,” in ICCV, 2015.
- [3] S. Ji, W. Xu, M. Yang, and K. Yu. “3d convolutional neural networks for human action recognition,” IEEE Trans. Pattern Analysis and Machine Intelligence, 2013.
- [4] C. Feichtenhofer, A. Pinz, and R. P. Wildes. “Spatiotemporal multiplier networks for video action recognition,” in 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 2017, pp. 7445–7454.
- [5] K. Simonyan and A. Zisserman. “two-stream convolutional networks for action recognition in videos,” in NIPS, 2014.
- [6] G. Varol, I. Laptev, and C. Schmid. “Long-term temporal convolutions for action recognition,” IEEE Transactions on Pattern Analysis and Machine Intelligence, 2017.
- [7] Zhe Cao and Gines Hidalgo and Tomas Simon and Shih-En Wei and Yaser Sheikh. “OpenPose: realtime multi-person 2D pose estimation using Part Affinity Fields,” arXiv preprint arXiv:1812.08008, 2018.
- [8] Zhe Cao and Tomas Simon and Shih-En Wei and Yaser Sheikh. “Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields,” CVPR, 2017.
- [9] Tomas Simon and Hanbyul Joo and Iain Matthews and Yaser Sheikh. “Hand Keypoint Detection in Single Images using Multiview Bootstrapping,” CVPR, 2017.
- [10] Shih-En Wei and Varun Ramakrishna and Takeo Kanade and Yaser Sheikh. “Convolutional pose machines,” CVPR, 2016.