

# 電腦網路導論期末報告

劉知穎 B07901039

陳韋丞 B07901060

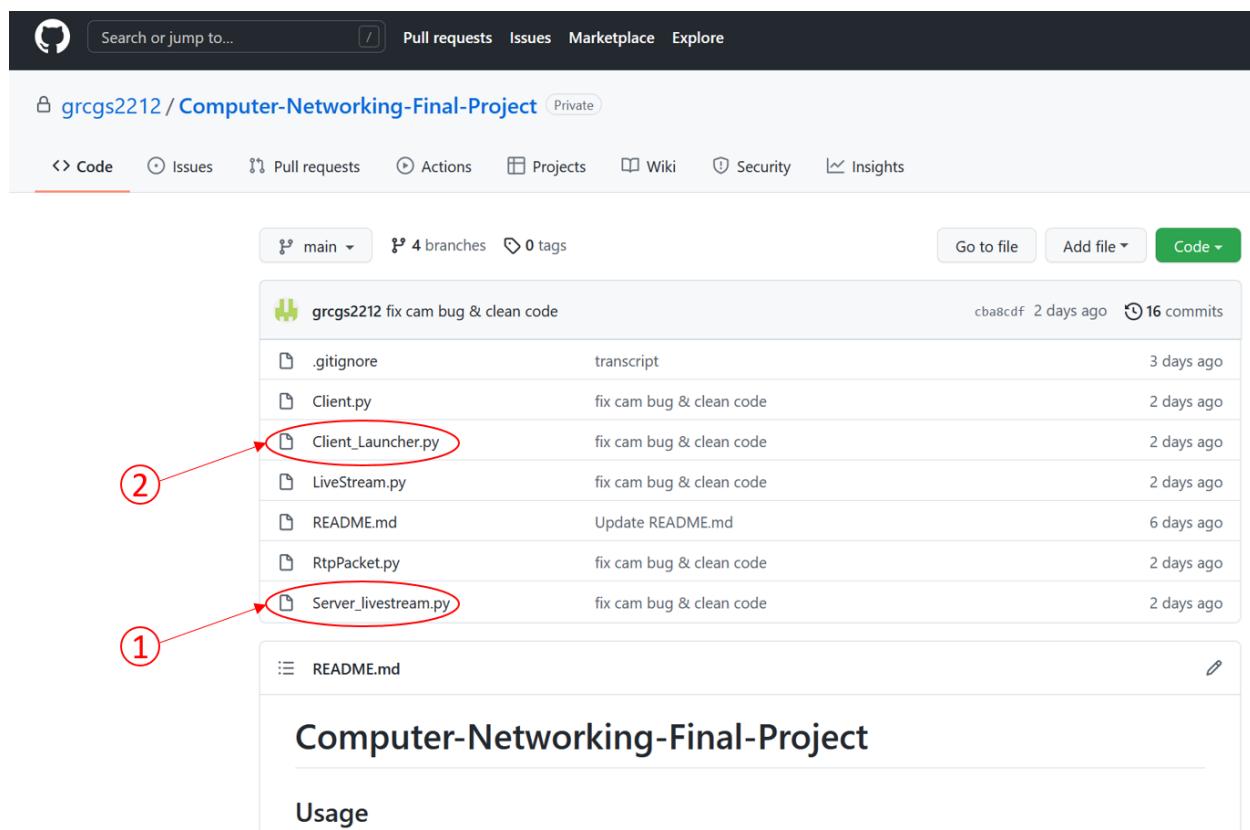
廖甜雅 B07901113

## 簡介

本次期末專題我們做了一個視訊通話的程式。起初我們先根據期末規定寫了一個video streaming的程式，並增加了回放以及進度條的功能，成功後，我們將其進一步改進成live streaming並增加了聲音傳輸；如此，我們便打下了視訊通話的基礎，其後我們為了解決通訊品質不穩定的問題，又增加了開啟字幕的功能，經過一番測試後，最終成功實現了兩人的視訊通話。惟通訊延遲平均大概2.5秒，仍有改進的空間；但大體而言溝通無礙。詳請請見我們的demo影片。

 Computer Networking Final Project Demo

## 使用介紹



The screenshot shows a GitHub repository page for 'grcgs2212 / Computer-Networking-Final-Project'. The 'Code' tab is selected. The main area displays a list of files in the 'main' branch:

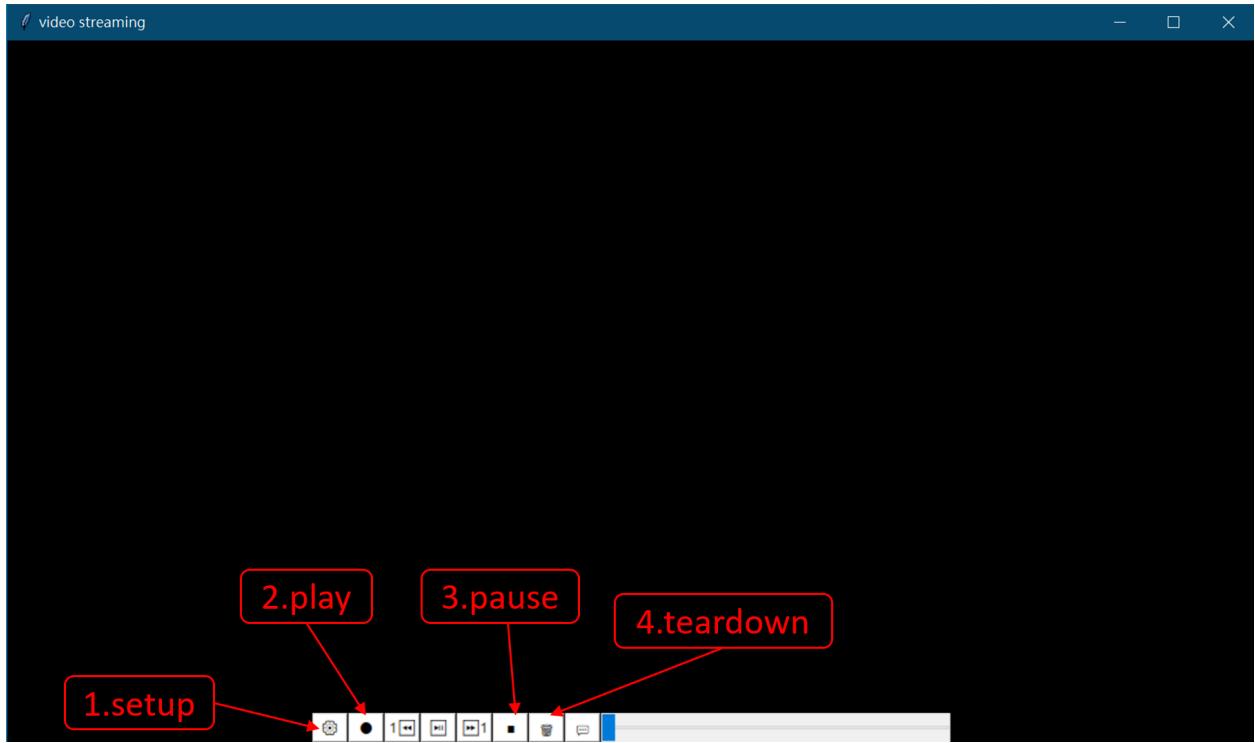
File	Description	Last Commit
.gitignore	transcript	3 days ago
Client.py	fix cam bug & clean code	2 days ago
Client_Launcher.py	fix cam bug & clean code	2 days ago
LiveStream.py	fix cam bug & clean code	2 days ago
README.md	Update README.md	6 days ago
RtpPacket.py	fix cam bug & clean code	2 days ago
Server_livestream.py	fix cam bug & clean code	2 days ago

Two specific files are highlighted with red circles and numbered: 'Client\_Launcher.py' (circled with ②) and 'Server\_livestream.py' (circled with ①). Below the file list, there is a preview of the 'README.md' file content:

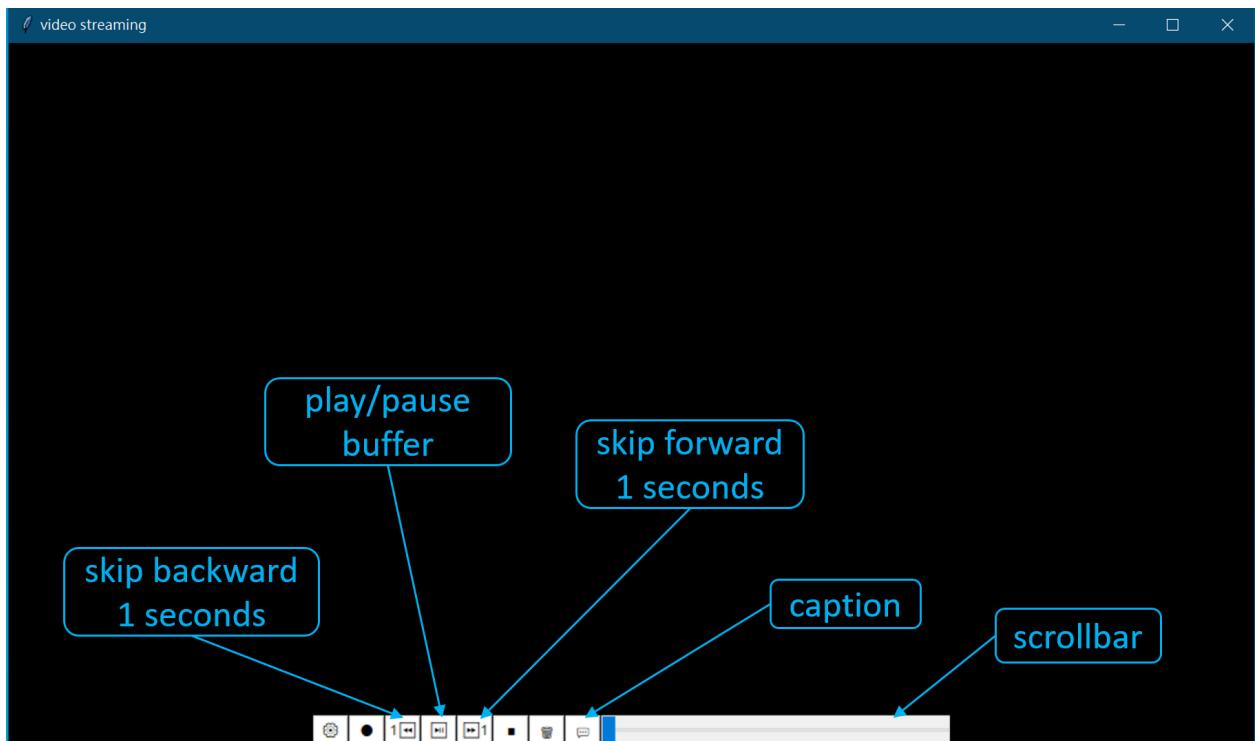
Computer-Networking-Final-Project

Usage

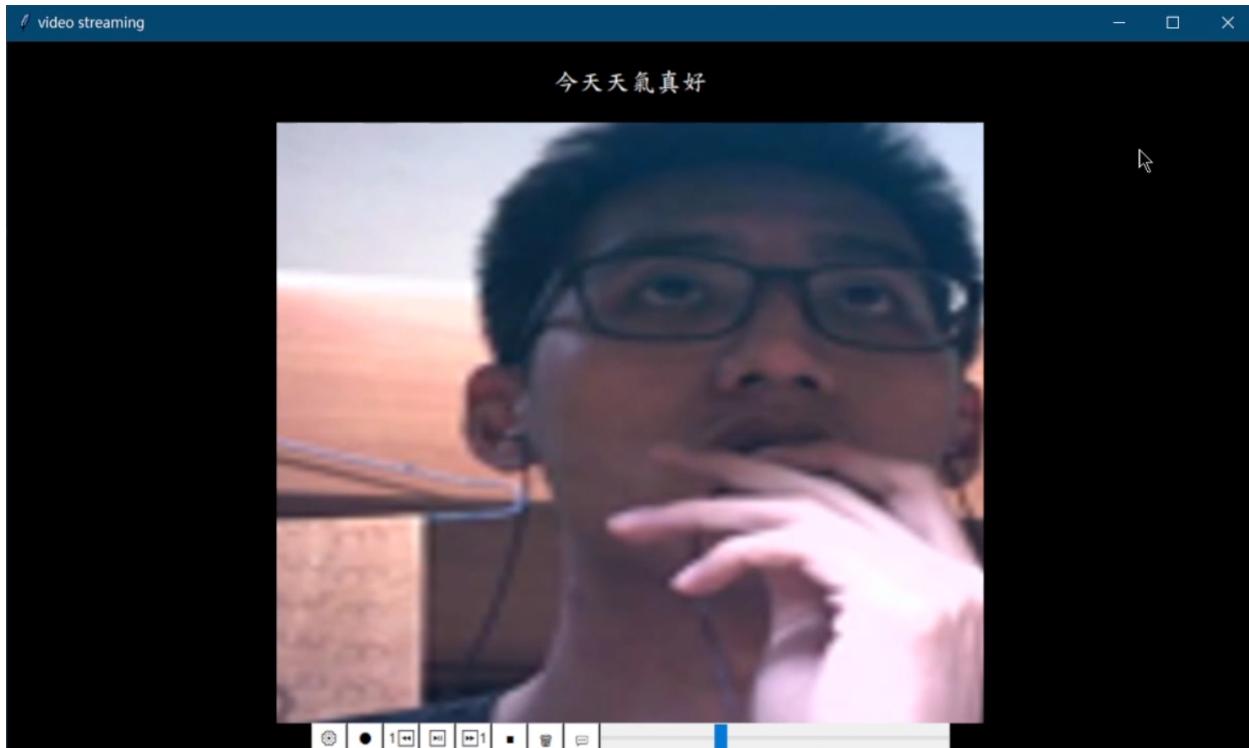
本次視訊通話程式在Python環境下運行，會使用到numpy、opencv-python、PyAudio、以及SpeechRecognition等模組。安裝完畢之後，播放者先執行Server\_livestream.py檔。待terminal顯示RTSP socket listening...之後，接收者再執行Client\_Launcher.py檔，便會跳出如下視窗。



根據RTSP protocol, client可以向server傳送setup、play、pause、teardown四種指令，我們以不同的符號來表示這些功能。另外，為了達成直播回放的功能，我們設計先將server回傳的資料儲存在buffer之中，而不直接播放。因此，協定中的play更趨像是錄製功能而非播放功能，故以錄製符號表示；而pause則更趨近於停止錄製功能，故以停止符號表示。



如上所述，我們設計先將直播資料錄製於buffer儲存，所以接收者需點擊上圖中間play/pause按鈕來播放或暫停影片。我們希望接收者可以有很高的自由度，故設計了快進按鈕、倒退按鈕、以及進度條(scrollbar)讓使用者可以任意調整自己想看的直播片段。另外，因為通訊品質會受限於網路環境，有時可能會聽不清楚對方在說甚麼，因此我們還設計了字幕功能。點擊caption符號開啟字幕功能後，server端會辨識播放者的聲音並轉換成文字，然後傳送給接收者。



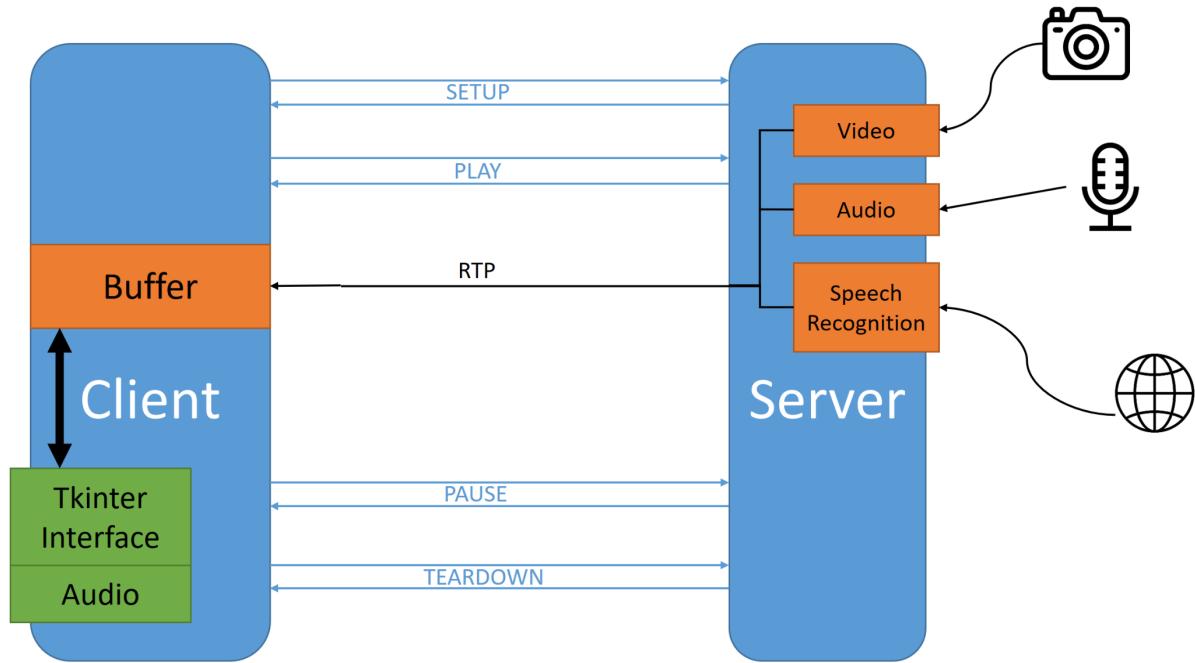
如圖為實際直播畫面。可以注意到，若我們可以另外開通一個logical channel，並交換播放者與接收者，便能實現視訊通話的功能。操作方法僅需兩方交換，其餘不變。



上圖為我們實際視訊通話的成果，平均延遲大概2.5秒左右。詳情請見我們的demo影片。

## 程式介紹

在這個程式中，直播資料的提供者為Server端，直播資料的接收者為Client端。兩者以TCP傳遞請求與回應訊息，使用RTSP protocol。在播放直播的過程中，Server端的直播畫面、聲音、語音辨識字幕是以RTP傳送給Client端，而Client端處理從Server端接收的直播資料，將畫面顯示於以tkinter建立的使用者介面，同時播出音訊。



### <RTP packet>

```
1 import sys
2
3 HEADER_SIZE = 2 # only consider 2 bytes seq num
4
5 class RtpPacket:
6     def __init__(self):
7         self.seqnum = None
8         self.payload = None
9
10    def encode(self, seqnum, payload):
11        self.seqnum = bytearray(seqnum.to_bytes(HEADER_SIZE, sys.byteorder))
12        self.payload = payload
13
14    def getPacket(self):
15        return self.seqnum + self.payload
16
17    def decode(self, bytestring):
18        self.seqnum = bytestring[:HEADER_SIZE]
19        self.payload = bytestring[HEADER_SIZE:]
20
21    def getSeqNum(self):
22        return int.from_bytes(self.seqnum, sys.byteorder)
23
24    def getPayload(self):
25        return self.payload
```

RtpPacket.py為Server端與Client端皆會使用到的檔案，檔案中定義class RtpPacket，包含RTP packet結構 -- header(序列號)+payload，以及encode、decode的方法。

### <Server>

Server端使用的Python程式碼包含Server\_livestream.py、RtpPacket.py與LiveStream.py。其中Server\_livestream.py為主程式；RtpPacket.py如前面所述定義RtpPacket；LiveStream.py定義LiveStreamVideo與LiveStreamAudio兩個class。以下是LiveStream.py的內容與解釋。

```
1 import cv2
2 import pyaudio
3 import time
4 from PIL import Image
5 import numpy as np
6
7
8 class LiveStreamVideo:
9     def __init__(self):
10         self.framNum = 0
11         self.cap = cv2.VideoCapture(0)
12
13     def getNextFrame(self):
14         self.framNum += 1
15         ret, frame = self.cap.read()
16         try:
17             frame = Image.fromarray(frame)
18             frame = frame.resize((160, 120))
19             frame = np.array(frame)
20             frame = frame.tobytes()
21             print('----- Next Video Frame (# {}), length: {} -----'.format(self.framNum, len(frame)))
22         except: pass
23         return frame
```

LiveStreamVideo負責接收設備攝影機所讀取的畫面。初始化時會開啟攝影機鏡頭；方法getNextFrame()會回傳當下視訊畫面的byte array，為避免攝影機鏡頭畫質過高導致訊息過長無法順利傳送，在讀取畫面後會先轉換成160\*120的大小再加以編碼回傳。經過測試，視訊畫面為每秒30幀。

```

25  class LiveStreamAudio:
26      def __init__(self):
27          self.framNum = 0
28          self.NUM_FRAME_PER_CHUNK = 10
29          self.RATE = 44100 # num frames per second
30          self.CHUNK = int(44100 / 30 * self.NUM_FRAME_PER_CHUNK) # num frame
31          self.FORMAT = pyaudio.paInt16
32          self.CHANNELS = 2
33          self.buffer = []
34          self.p = pyaudio.PyAudio()
35          self.stream = self.p.open(format=self.FORMAT,
36              channels=self.CHANNELS,
37              rate=self.RATE,
38              input=True,
39              frames_per_buffer=self.CHUNK)
40          # self.stream.stop_stream() # for the first chunk delay
41
42      def getNextChunk(self):
43          # self.stream.start_stream()
44          data = self.stream.read(self.CHUNK)
45          self.buffer.append(data)
46          # self.stream.stop_stream()
47          self.framNum += self.NUM_FRAME_PER_CHUNK
48          print('\t\t\t---- Next Audio Chunk (# {}), length: {} ----'.format(self.framNum, len(data)))
49          return data

```

LiveStreamAudio負責接收麥克風所讀取的聲音。初始化時會開始接收麥克風聲音；方法getNextChunk()會回傳self.CHUNK大小的一段音訊。由於音訊無法以太短的單位接收，否則會導致Client端的聲音聽起來斷斷續續，大幅降低音質，因此每次傳送約10幀視訊長度的音訊，由此計算 self.CHUNK = (每幀視訊時間與每幀音訊時間的比值) \* 10 = (44100 / 30) \* 10，取整數。

以下為主程式Server\_livestream.py的內容。為便於理解，下圖先展示Server\_livestream.py 開始執行時首先執行的內容，再解釋其餘部分的程式碼。

```

144  if __name__ == '__main__':
145
146      HOST, PORT = '127.0.0.1', 8888 # For local network (server and client seperate on 2 computers)
147      # HOST, PORT = '127.0.0.1', 8888 # For one host (on the same computer)
148      rtsp_socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM) # RTSP: TCP socket
149      rtsp_socket.bind((HOST, PORT))
150
151      # to avoid delay, open live stream before listening
152      live_stream_video = LiveStreamVideo()
153      live_stream_audio = LiveStreamAudio()
154      if not os.path.exists('./cache'):
155          os.makedirs('./cache')
156      print('RTSP socket listening...')
157      rtsp_socket.listen(5)
158      while True:
159          rtsp_client, addr = rtsp_socket.accept() # this accept {SockID,tuple object},tuple object = {clinet_addr,intNum}!!!
160          ServerWorker(rtsp_client, addr, live_stream_video, live_stream_audio).run()

```

在剛開始執行主程式的時候，建立Server與Client對話的 RTSP socket，並開始等待Client的連線。為避免視訊延遲，在開始listen之前先開啟鏡頭畫面與麥克風聲音，並創建一個cache資料夾以便之後暫存欲進行語音辨識的聲音資料。

```

1 import os
2 import sys
3 import socket
4 import threading
5 from random import randint
6 from LiveStream import LiveStreamVideo, LiveStreamAudio
7 from RtpPacket import RtpPacket
8 import speech_recognition as sr
9 import wave
10 import time
11
12
13 class ServerWorker:
14     def __init__(self, socket, clientAddr, live_stream_video, live_stream_audio):
15         self.rtsp_socket = socket
16         self.rtp_socket_video = None
17         self.rtp_socket_audio = None
18         self.rtp_socket_word = None
19         self.rtp_addr = clientAddr[0]
20         self.rtp_port_video = None
21         self.rtp_port_audio = None
22         self.rtp_port_word = None
23         self.state = 'INIT' # 'INIT', 'READY', 'PLAYING'
24         self.live_stream_video = live_stream_video
25         self.live_stream_audio = live_stream_audio
26         self.session = None
27         self.event = None
28         self.worker_video = None
29         self.worker_audio = None
30         self.framNum_video = 0 # for aligning video and audio
31         self.framNum_audio = 0
32
33     def run(self):
34         threading.Thread(target=self.receiveRTSPRequest).start()
35
36     def receiveRTSPRequest(self):
37         while True:
38             data = self.rtsp_socket.recv(1024)
39             if data:
40                 self.processRTSPRequest(data)

```

與Client成功連線時，初始化所有參數，起始狀態為INIT，並開啟一個新的執行緒聆聽Client的請求。當收到請求時，呼叫processRTSPRequest()以處理請求。

```

42     def processRTSPRequest(self, data):
43         print('### RTSP request received: {}'.format(data))
44         request = data.decode().split('\n')
45         requestType = request[0].split(' ')[0]
46         seqNum = int(request[1])
47         if requestType == 'SETUP':
48             if self.state == 'INIT':
49                 print('SETUP request received')
50                 self.state = 'READY'
51                 self.session = randint(100000, 999999)
52                 self.rtp_port_video = int(request[2].split(' ')[-3])
53                 self.rtp_port_audio = int(request[2].split(' ')[-2])
54                 self.rtp_port_word = int(request[2].split(' ')[-1])
55                 self.replyRTSP('OK_200', seqNum)
56         elif requestType == 'PLAY':
57             if self.state == 'READY':
58                 print('PLAY request received')
59                 self.state = 'PLAYING'
60                 self.rtp_socket_video = socket.socket(socket.AF_INET, socket.SOCK_DGRAM)
61                 self.rtp_socket_audio = socket.socket(socket.AF_INET, socket.SOCK_DGRAM)
62                 self.rtp_socket_word = socket.socket(socket.AF_INET, socket.SOCK_DGRAM)
63                 self.event = threading.Event()
64                 self.worker = threading.Thread(target=self.sendRTP_video_and_audio)
65                 self.worker.start()
66                 self.replyRTSP('OK_200', seqNum)
67         elif requestType == 'PAUSE':
68             if self.state == 'PLAYING':
69                 print('PAUSE request received')
70                 self.state = 'READY'
71                 self.event.set()
72                 self.replyRTSP('OK_200', seqNum)
73         elif requestType == 'TEARDOWN':
74             print('TEARDOWN request received')
75             self.event.set()
76             self.replyRTSP('OK_200', seqNum)
77             self.worker.join()
78             self.rtp_socket_video.close()
79             self.rtp_socket_audio.close()
80             self.rtp_socket_word.close()
81         else: pass
82
138     def replyRTSP(self, code, seqNum):
139         if code == 'OK_200':
140             reply = 'RTSP/1.0 200 OK\nCSeq: {}\nSession: {}'.format(seqNum, self.session).encode()
141             self.rtsp_socket.send(reply)

```

Client可能傳送setup、play、pause、和teardown請求，完成請求後，傳送'OK\_200'至Client。

收到setup請求時，切換狀態到READY，並設定收到的視訊、音訊與字幕的rtp port。

收到play請求時，切換狀態到PLAYING，建立視訊、音訊與字幕的rtp socket，並開啟新的執行緒執行sendRTP\_video\_and\_audio()，持續傳送當前視訊、音訊以及字幕內容給Client。

收到pause請求時，回到READY狀態並暫停傳送影音與字幕。

收到teardown請求時，暫停傳送影音與字幕，待傳送資料的執行緒結束之後關閉所有socket。

```

83     def asr(self, framNum):
84         WAV_OUTPUT_FILE = os.path.join('cache', '{}.wav'.format(self.session))
85         wf = wave.open(WAV_OUTPUT_FILE, 'wb')
86         wf.setnchannels(self.live_stream_audio.CHANNELS)
87         wf.setsampwidth(self.live_stream_audio.p.get_sample_size(self.live_stream_audio.FORMAT))
88         wf.setframerate(self.live_stream_audio.RATE)
89         wf.writeframes(b''.join(self.live_stream_audio.buffer[framNum//10-10:framNum//10]))
90         wf.close()
91         r = sr.Recognizer()
92         WAV = sr.AudioFile(WAV_OUTPUT_FILE)
93         with WAV as source:
94             audio = r.record(source)
95             output = r.recognize_google(audio, language="zh-TW", show_all=True)
96             if len(output) > 0:
97                 sent = output['alternative'][0]['transcript']
98             else:
99                 sent = ""
100            self.rtp_socket_word.sendto(self.makeRtpPacket(sent.encode(), framNum), (self.rtp_addr, self.rtp_port_word))
101            # print(framNum, sent)
102

```

asr()對當前音訊執行語音辨識，並透過字幕port傳送字幕的rtp packet。

```

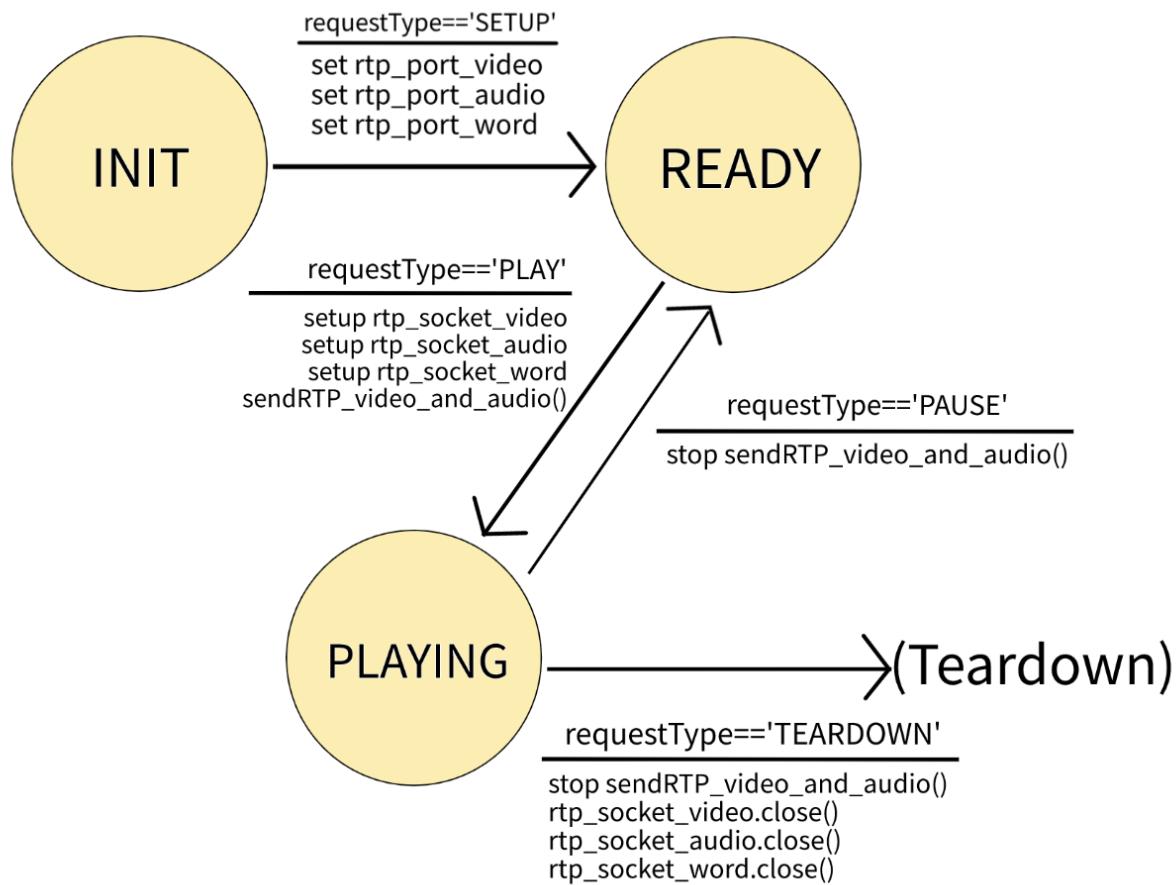
103     def sendRTP_video_and_audio(self):
104         while True:
105             tot_start_time = time.time()
106             if self.event.isSet():
107                 asr_thread.join()
108                 break
109             # audio, delay 0.3 sec
110             start_time=time.time()
111             data = self.live_stream_audio.getNextChunk() # delay 0-0.008 sec
112             framNum = self.live_stream_video.framNum
113             if framNum > 9:
114                 asr_thread = threading.Thread(target=self.asr, args=(framNum,))
115                 asr_thread.start()
116                 self.rtp_socket_audio.sendto(self.makeRtpPacket(data, framNum), (self.rtp_addr, self.rtp_port_audio))
117                 end_time = time.time()
118                 print('Total Audio Delay: {}'.format(end_time-start_time))
119             # video, delay 0.002
120             for i in range(10):
121                 if self.event.isSet():
122                     break
123                 start_time=time.time()
124                 data = self.live_stream_video.getNextFrame() # delay 0.002 sec
125                 framNum = self.live_stream_video.framNum
126                 self.framNum_video = framNum
127                 self.rtp_socket_video.sendto(self.makeRtpPacket(data, framNum), (self.rtp_addr, self.rtp_port_video))
128                 end_time = time.time()
129                 print('Total Video Delay: {}'.format(end_time-start_time))
130             tot_end_time = time.time()
131             # print('Total 10 frame delay: {} ({}).format(tot_end_time-tot_start_time, tot_end_time-tot_start_time-1/3)')
132
133     def makeRtpPacket(self, payload, framNum):
134         rtpPacket = RtpPacket()
135         rtpPacket.encode(framNum, payload)
136         return rtpPacket.getPacket()
137

```

當Server狀態在PLAYING時，會持續執行sendRTP\_video\_and\_audio()，直到pause或teardown請求暫停傳送。其中上圖的109-118行負責接收並傳送當前音訊，同時呼叫asr()辨識與傳送字幕；119-129行負責傳送視訊畫面。因為接收的音訊長度與接收的視訊長度比為10:1，所以在這裡以1段音訊:10幀視訊的頻率傳送聲音與畫面。

makeRtpPacket以framNum作為header，將聲音、畫面或字幕製作成 RTP packet。

下圖為Server的Finite state machine:



<Client>

```
1 v import sys
2   import tkinter
3   from client import Client
4
5
6 v if __name__ == '__main__':
7
8     client_rtp_port_video = 6666
9     client_rtp_port_audio = 7777
10    client_rtp_port_word = 5555
11    filename = "video.mjpeg" # for video streaming
12    root = tkinter.Tk()
13    root.geometry("1072x603")
14    root.configure(background='black')
15    app = Client(root, client_rtp_port_video, client_rtp_port_audio, client_rtp_port_word, filename)
16    root.mainloop()
17
```

如上圖, Client\_Launcher.py 會產生視訊、音訊、和字幕的rtp port , 並利用tkinter產生使用者介面。視訊通話不會用到filename。下面我們接著看Client.py檔。

```

1 import socket
2 import time
3 import threading
4 import numpy as np
5 from tkinter import*
6 from tkinter import messagebox as tkMessageBox
7 from tkinter import ttk
8 from RtpPacket import RtpPacket
9 from PIL import Image, ImageTk
10 import pyaudio
11
12
13 T_SEC = 1
14 BUFFER_SIZE = 1000
15 # SERVER_HOST, SERVER_PORT = '192.168.0.2', 8888 # communication on network
16 SERVER_HOST, SERVER_PORT = '127.0.0.1', 8888 # on the same host
17 CLIENT_HOST = '0.0.0.0'
18

```

如上圖, Client.py會產生五個constant , 其中T\_SEC表快進或倒退秒數。

```

19 class Client:
20     def __init__(self, master, rtpPort_video, rtpPort_audio, rtpPort_word, filename):
21         ## tkinter
22         self.master = master
23         self.master.title("video streaming")
24         self.master.protocol("WM_DELETE_WINDOW", self.handler) # handle closing GUI window
25         self.playMovieButton = False # toggle
26         self.tkwindow()
27         self.audio()
28         ## videotreamer
29         self.buffer_video = [] # buffer to realize replay function
30         self.buffer_audio = [] # buffer to realize replay function
31         self.buffer_word = [] # buffer to realize replay function
32         self.playIndex_video = 0
33         self.playIndex_audio = 0
34         self.playIndex_word = 0
35         self.buffNum = 0 # buffer will reset when full, buffNum = Nth buffer we are now at
36         self.SHOW_TRANSCRIPT = False
37         self.RESET_STATE = False
38         ## RTP
39         self.rtpPort_video = rtpPort_video
40         self.rtpPort_audio = rtpPort_audio
41         self.rtpPort_word = rtpPort_word
42         self.playRequestEvent = threading.Event()
43         self.runEvent = threading.Event()
44         ## RTSP
45         self.state = 'INIT' # 'INIT', 'READY', 'PLAYING'
46         self.sessionID = 0
47         self.filename = filename
48         self.requestSent = None
49         self.SETUP_STR = "SETUP {} \n1\n RTSP/1.0 RTP/UDP {} {} {}".format(self.filename, self.rtpPort_video, self.rtpPort_audio, self.rtpPort_word)
50         self.PLAY_STR = 'PLAY \n2'
51         self.PAUSE_STR = 'PAUSE \n3'
52         self.TEARDOWN_STR = 'TEARDOWN \n4'
53         self.connectToServer()
54         self.tearEvent = threading.Event()
55

```

Client端會同時開三個thread和buffer處理並儲存視訊、音訊、和字幕資料, 第32-34行的playIndex 紀錄螢幕上顯示的影片播放位置。

```
56  def connectToServer(self):
57      ## RTSP / TCP session
58      self.rtsp_socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
59      self.rtsp_socket.connect((SERVER_HOST, SERVER_PORT))
60
61  def openRTPsocket(self):
62      ## video
63      self.rtp_socket_video = socket.socket(socket.AF_INET,socket.SOCK_DGRAM)
64      self.rtp_socket_video.settimeout(0.5)
65      self.rtp_socket_video.bind((CLIENT_HOST, self.rtpPort_video))
66      ## audio
67      self.rtp_socket_audio = socket.socket(socket.AF_INET,socket.SOCK_DGRAM)
68      self.rtp_socket_audio.settimeout(0.5)
69      self.rtp_socket_audio.bind((CLIENT_HOST, self.rtpPort_audio))
70      ## word
71      self.rtp_socket_word = socket.socket(socket.AF_INET,socket.SOCK_DGRAM)
72      self.rtp_socket_word.settimeout(0.5)
73      self.rtp_socket_word.bind((CLIENT_HOST, self.rtpPort_word))
74
```

如上圖，初始時便會直接呼叫connectToServer()，其利用TCP與server端根據RTSP protocol通訊。setup時會呼叫openRTPsocket()，其會用UDP建立視訊、音訊、以及字幕的rtp\_socket。

```

75    def recvRtspReply(self):
76        while not self.tearEvent.isSet():
77            try:
78                data = self.rtsp_socket.recv(1024)
79                if data:
80                    print('### RTSP reply received: {}'.format(data))
81                    self.parseRtspReply(data)
82            except: continue # handle time out
83
84    def parseRtspReply(self, data):
85        lines = data.decode().split('\n')
86        session = int(lines[2].split(' ')[1])
87        if self.sessionID == 0: self.sessionID = session # initiate
88        if self.sessionID == session and int(lines[0].split(' ')[1]) == 200:
89            if self.requestSent == 'SETUP':
90                self.state = 'READY'
91                self.openRTPsocket()
92            elif self.requestSent == 'PLAY':
93                self.state = 'PLAYING'
94            elif self.requestSent == 'PAUSE':
95                self.state = 'READY'
96            elif self.requestSent == 'TEARDOWN':
97                self.state = 'INIT'
98                self.rtp_socket_video.shutdown(socket.SHUT_RDWR)
99                self.rtp_socket_video.close()
100               self.rtp_socket_audio.shutdown(socket.SHUT_RDWR)
101               self.rtp_socket_audio.close()
102               self.rtp_socket_word.shutdown(socket.SHUT_RDWR)
103               self.rtp_socket_word.close()
104            else: print("Session not in order")
105

```

如上圖，只要client不執行tear，recvRtspReply()便會一直接收RTSP資訊。當client接收到RTSP資訊後便會呼叫parseRtspReply()來處理接收資料，並改變client狀態。

```
106     def listenRtp_video(self):
107         while not (self.playRequestEvent.isSet() or self.tearEvent.isSet()):
108             try:
109                 data, addr = self.rtp_socket_video.recvfrom(80000)
110                 if data:
111                     rtpPacket = RtpPacket()
112                     rtpPacket.decode(data)
113                     currFrameNbr = rtpPacket.getSeqNum()
114                     ## late packet
115                     if currFrameNbr < self.buffNum * BUFFER_SIZE + len(self.buffer_video):
116                         continue
117                     ## correct packet
118                     else:
119                         self.buffer_video.append(rtpPacket.getPayload())
120                         if len(self.buffer_video) > BUFFER_SIZE:
121                             self.reset()
122             except: continue # handle out of time
```

如上圖，只要沒有pause或tear，listenRtp\_video()便會不斷接收server端傳來的視訊資料，並儲存於buffer之中。注意到第120行，BUFFER\_SIZE只要記憶體允許的話可以開得非常大，亦即可以將整段直播完全錄下來。但因為我們還沒找到非常有效得對齊視訊與音訊的方式，所以若不讓buffer reset的話，誤差會不斷累積；根據我們的實測，35秒的通話大概會造承3到4秒視訊音訊誤差。

```

124     def listenRtp_word(self):
125         while not (self.playRequestEvent.isSet() or self.tearEvent.isSet()):
126             try:
127                 data, addr = self.rtp_socket_word.recvfrom(4096)
128                 if data:
129                     rtpPacket = RtpPacket()
130                     rtpPacket.decode(data)
131                     currFrameNbr = rtpPacket.getSeqNum()
132                     ## late packet
133                     if currFrameNbr < self.buffNum * BUFFER_SIZE + len(self.buffer_word):
134                         continue
135                     ## correct packet
136                     else:
137                         self.buffer_word.append(rtpPacket.getPayload())
138                         if len(self.buffer_word) > BUFFER_SIZE//10:
139                             self.reset()
140             except: continue
141
142     def listenRtp_audio(self):
143         while not (self.playRequestEvent.isSet() or self.tearEvent.isSet()):
144             try:
145                 data, addr = self.rtp_socket_audio.recvfrom(80000)
146                 if data:
147                     rtpPacket = RtpPacket()
148                     rtpPacket.decode(data)
149                     currFrameNbr = rtpPacket.getSeqNum()
150                     ## late packet
151                     if currFrameNbr < self.buffNum * BUFFER_SIZE + len(self.buffer_audio):
152                         continue
153                     ## correct packet
154                     else:
155                         self.buffer_audio.append(rtpPacket.getPayload())
156                         if len(self.buffer_audio) > BUFFER_SIZE//10:
157                             self.reset()
158             except: continue # handle out of time
159

```

如上圖, listenRtp\_word()以及listenRtp\_audio()的功能與前述listenRtp\_video()非常類似, 惟注意到138和156行, 可以發現BUFFER\_SIZE//10與視訊的reset條件不同。原因是server傳過來的視訊和音訊的每一個chunk的時間比為1:10。

```

160     def reset(self):
161         ## clean both buffer and update buffer number (Nth)
162         self.buffer_video = []
163         self.buffer_audio = []
164         self.buffer_word = []
165         self.playIndex_video = 0
166         self.playIndex_audio = 0
167         self.playIndex_word = 0
168         self.bar.set(self.playIndex_video)
169         self.buffNum += 1
170

```

如上圖, reset()將buffer清空並使 index歸零。self.bar為使用者介面之進度條。

```

171     def setupRequest(self):
172         if self.state == 'INIT':
173             threading.Thread(target=self.recvRtspReply, daemon=True).start()
174             self.rtsp_socket.send(self.SETUP_STR.encode())
175             self.requestSent = 'SETUP'
176
177     def playRequest(self):
178         if self.state == 'READY':
179             self.playRequestEvent.clear()
180             threading.Thread(target=self.listenRtp_video, daemon=True).start()
181             threading.Thread(target=self.listenRtp_audio, daemon=True).start()
182             threading.Thread(target=self.listenRtp_word, daemon=True).start()
183             self.rtsp_socket.send(self.PLAY_STR.encode())
184             self.requestSent = 'PLAY'
185
186     def pauseRequest(self):
187         if self.state == 'PLAYING':
188             self.playRequestEvent.set()
189             self.rtsp_socket.send(self.PAUSE_STR.encode())
190             self.requestSent = 'PAUSE'
191
192     def tearRequest(self):
193         self.tearEvent.set()
194         self.rtsp_socket.send(self.TEARDOWN_STR.encode())
195         self.requestSent = 'TEARDOWN'
196         self.master.destroy() # Close the gui window
197         self.master.quit()
198

```

上圖的四個函式分別處理RTSP中的setup、play、pause、和teardown請求。執行pauseRequest()和tearRequest()時，會分別將self.playRequestEvent和self.tearEvent設成True，讓前述listenRtp\_()或recvRtspReply()迴圈停止。注意到這邊無法使影片播放，而僅會進行資料處理並存進buffer。

```

199     def playMovie(self):
200         self.runEvent.clear()
201         self.playIndex_video = ((len(self.buffer_video) - 1) // 10) * 10
202         self.playIndex_audio = self.playIndex_video + 10
203         self.playIndex_word = self.playIndex_audio
204         threading.Thread(target=self.run_video, daemon=True).start()
205         threading.Thread(target=self.run_audio, daemon=True).start()
206         threading.Thread(target=self.run_word, daemon=True).start()
207
208     def stopMovie(self):
209         self.runEvent.set()
210
211     def backwardMovie(self):
212         self.playIndex_video = ((self.playIndex_video - 30*T_SEC)//10) * 10
213         self.playIndex_audio = ((self.playIndex_video - 30*T_SEC)//10) * 10 + 10
214         self.playIndex_word = ((self.playIndex_video - 30*T_SEC)//10) * 10 + 10
215         if self.playIndex_video < 0 : self.playIndex_video = 0
216         if self.playIndex_audio < 0 : self.playIndex_audio = 0
217         if self.playIndex_word < 0 : self.playIndex_word = 0
218         self.bar.set(self.playIndex_video)
219         print("replay, frame # = ", self.playIndex_video)
220
221     def forwardMovie(self):
222         self.playIndex_video = ((self.playIndex_video + 30*T_SEC)//10) * 10
223         self.playIndex_audio = ((self.playIndex_video + 30*T_SEC)//10) * 10 + 10
224         self.playIndex_word = ((self.playIndex_video + 30*T_SEC)//10) * 10 + 10
225         if self.playIndex_video > len(self.buffer_video) : self.playIndex_video = len(self.buffer_video)
226         if self.playIndex_audio > len(self.buffer_audio) * 10 : self.playIndex_audio = len(self.buffer_audio) * 10
227         if self.playIndex_word > len(self.buffer_word) * 10 : self.playIndex_word = len(self.buffer_word) * 10
228         self.bar.set(self.playIndex_video)
229         print("replay, frame # = ", self.playIndex_video)
230
231     def play_pause_Movie(self):
232         if self.playMovieButton == True: # pause
233             self.playMovieButton = False
234             self.stopMovie()
235         else: # play
236             self.playMovieButton = True
237             self.playMovie()

```

如上圖，這邊實際進入控制影片播放的函式，可以注意到第200和209行的self.runEvent，其可以用來控制是否要進入播放影片的迴圈。注意到 $30*T\_SEC$ 的30為影片每秒貞數。

```

238     def run_word(self):
239         while not (self.runEvent.isSet() or self.tearEvent.isSet()):
240             if self.playIndex_word > len(self.buffer_word)*10:
241                 self.playIndex_word = len(self.buffer_word)*10
242             if self.playIndex_word < len(self.buffer_word)*10 and len(self.buffer_word) != 0:
243                 try:
244                     if self.SHOW_TRANSCRIPT:
245                         self.word.set(self.buffer_word[self.playIndex_word//10].decode())
246                         self.playIndex_audio += 10 # NUM_FRAME_PER_CHUNK
247                 except: pass
248                 time.sleep(1/30) # wait for listenRTP_audio thread to finish pushing
249                 self.playIndex_word = (int(self.bar.get())//10)*10 + 10 # update index by scrollbar
250
251     def run_audio(self):
252         while not (self.runEvent.isSet() or self.tearEvent.isSet()):
253             if self.playIndex_audio > len(self.buffer_audio)*10: self.playIndex_audio = len(self.buffer_audio)*10
254             if self.playIndex_audio < len(self.buffer_audio)*10 and len(self.buffer_audio) != 0:
255                 time.sleep(1/30) # wait for listenRTP_audio thread to finish pushing
256                 try:
257                     self.stream.write(self.buffer_audio[self.playIndex_audio//10])
258                 except: pass
259                 self.playIndex_audio += 10 # NUM FRAME PER CHUNK
260                 self.playIndex_audio = (int(self.bar.get())//10)*10 + 10 # update index by scrollbar
261
262     def run_video(self):
263         while not (self.runEvent.isSet() or self.tearEvent.isSet()):
264             self.playIndex_video = int(self.bar.get()) # update index by scrollbar
265             if self.playIndex_video > len(self.buffer_video):
266                 self.playIndex_video = len(self.buffer_video)
267                 self.bar.set(self.playIndex_video) # update scrollbar
268             if self.playIndex_video < len(self.buffer_video) and len(self.buffer_video) != 0:
269                 inp = np.asarray(bytarray(self.buffer_video[self.playIndex_video]), dtype=np.uint8).reshape(120, 160, 3)
270                 inp = inp[:, :, [2, 1, 0]]
271                 pilImage = Image.fromarray(inp)
272                 w = self.master.winfo_height() # current window height
273                 h = self.master.winfo_width() # current window width
274                 pilImage = pilImage.resize((h, w)) # rescale # , Image.ANTIALIAS
275                 imgtk = ImageTk.PhotoImage(pilImage)
276                 self.label.configure(image = imgtk, height=h, width=w)
277                 self.label.image = imgtk
278                 self.playIndex_video += 2
279                 self.bar.set(self.playIndex_video) # update scrollbar
280                 time.sleep(1/30) # 30 frames per second
281

```

執行playMovie()之後，會開三個thread分別執行上圖run\_word()、run\_audio()、以及run\_video()。只要不要點擊暫停播放或逕將socket teardown，三者便會不斷的執行。可以注意到run\_video()第272到274行，目的為調整播放影片的邊長以隨著視窗大小變化。再注意到第278行的playIndex，其每次加二而非一。原因是經過我們的實測結果，run\_video()跑一次迴圈需時0.06秒而run\_audio()跑一次迴圈需時0.3秒，視訊速度為音訊的20倍，但前述已經處理過10倍的問題，所以影格播放的速度僅需變兩倍。

```

289 def tkwindow(self):
290     """Build GUI."""
291     bottom_frame = Frame(self.master)
292     bottom_frame.pack(side=BOTTOM)
293     top_frame = Frame(self.master)
294     top_frame.pack(side=TOP)
295     # Create Setup button
296     self.setup = Button(bottom_frame, width=3, height=1, bg='white')
297     self.setup["text"] = "◎"
298     self.setup["command"] = self.setupRequest
299     self.setup.pack(side=LEFT)
300     # Create Play button
301     self.start = Button(bottom_frame, width=3, height=1, bg='white')
302     self.start["text"] = "●"
303     self.start["command"] = self.playRequest
304     self.start.pack(side=LEFT)
305     # Create backward button
306     self.backward = Button(bottom_frame, width=3, height=1, bg='white')
307     self.backward["text"] = "{} {}".format(T_SEC)
308     self.backward["command"] = self.backwardMovie
309     self.backward.pack(side=LEFT)
310     # Create playMovie button
311     self.play = Button(bottom_frame, width=3, height=1, bg='white')
312     self.play["text"] = "▣"
313     self.play["command"] = self.play_pause_Movie
314     self.play.pack(side=LEFT)
315     # Create forward button
316     self.forward = Button(bottom_frame, width=3, height=1, bg='white')
317     self.forward["text"] = "▣ {}".format(T_SEC)
318     self.forward["command"] = self.forwardMovie
319     self.forward.pack(side=LEFT)
320     # Create Pause button
321     self.pause = Button(bottom_frame, width=3, height=1, bg='white')
322     self.pause["text"] = "■"
323     self.pause["command"] = self.pauseRequest
324     self.pause.pack(side=LEFT)
325     # Create Teardown button
326     self.teardown = Button(bottom_frame, width=3, height=1, bg='white')
327     self.teardown["text"] = "☒"
328     self.teardown["command"] = self.tearRequest
329     self.teardown.pack(side=LEFT)
330     # Create Transcript button
331     self.teardown = Button(bottom_frame, width=3, height=1, bg='white')
332     self.teardown["text"] = "▣"
333     self.teardown["command"] = self.setTranscript
334     self.teardown.pack(side=LEFT)
335     # Create a bar
336     self.bar = ttk.Scale(bottom_frame, from_=0, to=BUFFER_SIZE, length=300)
337     self.bar.pack(side=LEFT)
338     # Create a text box for transcript
339     self.word = StringVar()
340     self.word.set("")
341     self.text = Label(top_frame, textvariable=self.word, bg='black', fg='white', width=20, height=3, font=(('標楷體', 16)))
342     self.text.pack(side=LEFT)
343     # Create a label to display the movie
344     self.label = Label(bg='black')
345     self.label.pack(side=TOP)

```

如上圖，利用tkinter模組，我們寫了一個使用者介面。

```

347 def handler(self): # handle closing GUI window directly without teardown button
348     """Handler on explicitly closing the GUI window."""
349     self.pauseRequest()
350     if tkMessageBox.askokcancel("Quit?", "Are you sure you want to quit?"):
351         self.tearRequest()
352     else: # When the user presses cancel, resume playing.
353         self.playRequest()

```

最後，handler會處理使用者按下視窗右上角叉叉按鈕時的情況。

## 問題與解決方式

### 1. 聲音影像不同步的問題

由於是livestream，我們的傳送的聲音是從麥克風錄製得來的，如果一次錄製的時長設計過短，聲音的品質會非常差，太長則無法用socket傳遞，且遇到有loss與delay的環境，就會一次損失一大段聲音。經過權衡，我們影片以30 fps傳遞，而聲音與影像傳遞頻率比率為1:10，以此設計相對應一個聲音chunk的長度。

然而，我們發現由於套件本身的delay，在server端當我們用兩個thread各自傳遞聲音與影像時，兩者的對其狀況並不如我們所想。在client端也是，由於影像要多做些轉換和tkinter的delay，影像遠比聲音慢得多。

為了解決這個問題，我們詳細研究每個套件的運作方式和delay。例如，我們發現pyaudio的stream.read()，會在背景持續填充buffer，因此這個函式本身的delay並不是我們預期的1/3秒，而是將盡0s。因此，在server端，我們可以用單一線程去輪流處理聲音與影像的傳送，再配合加入相應的delay，如此便能保證他們的比例為1:10，且影響速率為30 fps。

在client端，由於用單一線程依序處理聲音與影像的話，無法避免會聽到麥克風開開關關的聲音，因此我們決定以兩個線程各自去播影片與聲音。經過delay觀察，我們發現讓影像隔著播，就可以大約讓影像對齊。剩下偏差的部分，我們則在每次client端buffer滿了要清空時，強制讓聲音影像對齊。

總而言之，我們依靠設計multithread和調整delay，成功讓聲音與影像在server與client端對齊。在使用local host，不考慮網路延遲的情況下，影像與聲音的delay在1秒之內，且聲音的撥放很穩定。

## 加分項目

### 1. Livestream

Server傳送從電腦鏡頭拍攝的影像與從麥克風錄製的聲音到client。更進一步，只要兩端client互相連結到對方的server，就可以達到視訊通話。用localhost連結時delay在1秒之內，用網路線連結時則約2.5秒。

### 2. 使用者介面

使用tkinter刻劃精美的使用者介面。使用者可以用按鈕發送RTSP Request、影片撥放/暫停/快進/倒退、開啟或關閉字幕，也可使用進度條控制撥放時間點。

### 3. 控制影片撥放的功能

包括撥放/暫停/快進/倒退，使用進度條選擇撥放的時間點

4. 支援多個client連結到同個server
5. 字幕功能

使用speech\_recognition套件，在server端製作字幕。用額外一個UDP socket傳送。Client可以選擇開啟或關掉字幕。

## 未來展望

因為音訊不像視訊可以完整的切分成一貞或一個影格，故聽起來沒有到非常順暢，未來我們需要進一步改進我們的音訊傳輸品質。另外，我們或許可以壓縮資料的傳輸量，並在client端利用機器學習的方式將畫質提高，藉此改善平均2.5秒的通訊延遲。最後，我們的測試還可以更充足，希望可以進一步提高系統的穩定性並將所有的程式打包成執行檔，以利一般使用者使用。