

高雄科技大學

智慧商務系

AZURE 雲端服務認證與證照

第二次期中專題

影片字幕產生器

姓名：許智程

學號：C109152304

指導老師：謝文川 老師

目錄

目錄	i
1 研究動機與目的	1
2 專案規劃與設計	1
2.1 專案架構簡介	1
2.1.1 FFmpeg	1
2.1.2 Azure Cognitive Services	2
3 程式實作	2
3.1 抽取音檔	2
3.2 從音檔生成字幕	3
3.3 翻譯得到的字幕檔	7
3.4 把字幕跟影片結合	8
3.5 Log 與 Decorator	9
4 後記	10
參考資料	11

1 研究動機與目的

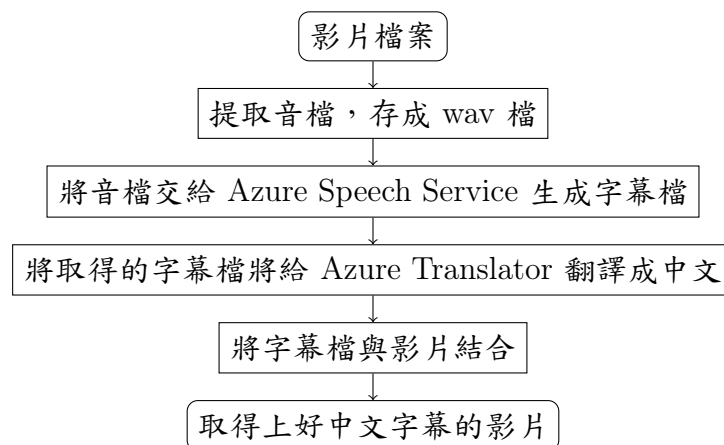
學習資訊技術時，常常得面臨查到的資料只有英文的部分，如果是文章的話當看不懂時還可以用 Chrome 的翻譯功能將文章翻譯成中文幫助理解，不過文章看久了總想看看別人實際下來實作，而在網路上多數的英文實作影片幾乎沒有中文字幕，還得面臨外國佬的口音很重的問題。

除了資訊技術的影片外，很多中文以外的內容通常也是蠻有料的，難道要因為語言的隔閡去放棄這些內容嗎？

為此，本次專題的目標就是利用在課堂上所學的 Azure AI 技巧來達成自動為影片生成字幕檔，並翻譯成中文，最後直接將字幕嵌到影片中。

2 專案規劃與設計

2.1 專案架構簡介



如上面的流程圖所示，首先會將影片的音檔提取出來，並存成 wav 檔，方便後續處理，接下來藉由 Azure Captioning 和 Azure Translation 來取得字幕並翻譯成中文，不過得到的結果會將字幕格式檔中紀錄時間的戳記中的標點符號全部由半形轉成全形，需要再經過處理，處理過後再將影片檔與字幕檔結合在一起就能取得上好中文字幕的影片了。

除此之外，還引入了 Log 的機制，紀錄了函數開始的時間與結束的時間，並且計算了耗費的秒數，下面將會簡單介紹所運用到的技術。

2.1.1 FFmpeg

FFmpeg 是一個開放原始碼的自由軟體，它包含了音訊和視訊多種格式的錄影、轉檔、串流功能，同時也是一個音訊與視訊格式轉換函式庫 (Library)，許多開源的工具都是基於 FFmpeg 打造的。FFmpeg 主要由以下幾個元件組成：

- CLI Application
 - ffmpeg：用來轉換不同格式的多媒體檔案，為 FFmpeg 最主要的工具
 - ffplay：基於 SDL(Simple DirectMedia Layer) 和 FFmpeg 函式庫的簡易播放器 [1]
 - ffprobe：簡易的多媒體串流分析器
- Libraries

- libavcodec：用來編碼與解碼不同的影音格式
- libavformat：用來處理承載影音串流的多媒體容器
- libavutil：一些程式上會用到的函式與資料結構
- libpostproc：對於視訊做前處理的函式庫
- libswscale：對於影像作縮放的函式庫
- libavfilter：對於影片與音訊做進一步處理，如倍速、水平翻轉、裁剪、加方框、疊加文字等功能
- libswresample：對於很多播放器，在輸出時會固定為一種格式（如 44100hz，雙聲道，16bit signed），因為多數設備能夠支持這些格式。這種情況下對於不同的多種輸入源，即需要進行音頻重採樣。而相對於舊式的 libavcodec，libswresample 能更輕鬆的處理。

運用這些程式與函數庫，我們可以方便地進行影音處理，例如說：從影片中提取音檔、將字幕嵌到影片中、影音轉檔..... 等都可以藉由 FFmpeg 完成，甚至可以用 FFmpeg 來剪片！

FFmpeg 如此的強大，可想而知的其複雜度也是不容小覷，於是不論是免費的還是付費的，許多人都對 FFmpeg 進行一定的封裝產生許多影音處理工具，連 python 都有輔助使用的 package 了呢！

在這次專題中，會藉由 ffmpeg-python[2] 這個在 python 中的輔助工具來進行取得影片的聲音檔與將字幕與影片檔壓在一起的動作。

2.1.2 Azure Cognitive Services

Azure Cognitive Services 是雲端式的 AI 服務，可使開發人員無須具備 AI 或 Data Science 的相關知識就可以將相關技術整合道專案中。這些服務可透過熱門開發語言（比如說：C#、Java、Python、JavaScript、golang..... 等）的 REST API 和 SDK 取得。Azure Cognitive Services 可讓開發人員使用認知解決方案，輕鬆地將認知功能新增至其應用程式，以進行查看、聆聽、說話和分析。

在這次專題中，會使用到 Azure Cognitive Services 中的 Speech Service 與 Translator 來達成生成影片的中文字幕的目的。

3 程式實作

3.1 抽取音檔

在 ffmpeg 中要拿到影片的音檔可說是非常容易，可以視作將影片轉成聲音檔，也就是說，輸入填上影片檔的路徑，輸出檔名只要是聲音檔的附檔名就可以了。為了方便進行下一步，我選擇轉成 wav 檔。如下面的程式所示：

```

1 from utils import Logger
2 import ffmpeg
3
4
5 @Logger
6 def splitwav(orginFilePath, targetFilePath):
7     (
8         ffmpeg.input(orginFilePath)
9         .output(targetFilePath)
10        .overwrite_output()

```

```
11     .run()
12 )
```

咦 @Logger 是甚麼？簡單來說，這是用一種叫做 Decorator 的技巧來讓紀錄程式執行的 Code 能抽成一個函數獨立出來，而不用黏在每一個函數之中。後面會再詳細說明。

3.2 從音檔生成字幕

很幸運的，在微軟的文檔 [3] 中有相關的範例，不過在看過範例程式碼後，發現他們是做成 CLI Tool 的形式，為了更好的與其他部分結合，需要動點手腳，範例 Code 節錄如下：

```
1 #
2 # Copyright (c) Microsoft. All rights reserved.
3 # Licensed under the MIT license. See LICENSE.md file in the project root
  for full license information.
4 #
5
6 # Notes:
7 # - Install the Speech SDK. Run:
8 # pip install azure-cognitiveservices-speech
9 # - The Python Speech SDK on Windows requires the Microsoft Visual C++
  Redistributable for Visual Studio 2015, 2017, 2019, or 2022 on the
  system. See:
10 # https://docs.microsoft.com/azure/cognitive-services/speech-service/
  quickstarts/setup-platform
11 # - Install gstreamer:
12 # https://docs.microsoft.com/azure/cognitive-services/speech-service/how-
  to-use-codec-compressed-audio-input-streams
13
14 # 顧及篇幅，只從範例程式全部4個檔案中節錄進入點的主要Class和實際執行的程式
15 # 如要查看全部的Code
16 # 請至: https://github.com/Azure-Samples/cognitive-services-speech-sdk/
  tree/master/scenarios/python/console/captioning
17 # file: captioning.py
18
19 from datetime import datetime, time, timezone
20 from itertools import groupby, pairwise
21 from os import linesep, remove
22 from os.path import exists
23 from pathlib import Path
24 from sys import argv
25 from time import sleep
26 from typing import Any, List, Optional
27 import wave
28 import azure.cognitiveservices.speech as speechsdk # type: ignore
29 import caption_helper
30 import helper
31 import user_config_helper
32
```

```

33 class Captioning(object) :
34     def __init__(self) :
35         self._user_config = user_config_helper.user_config_from_args(USAGE
36         )
37         self._srt_sequence_number = 1
38         self._previous_caption : Optional[caption_helper.Caption] = None
39         self._previous_end_time : Optional[time] = None
40         self._previous_result_is_recognized = False
41         self._recognized_lines : List[str] = []
42         self._offline_results : List[speechsdk.SpeechRecognitionResult] =
43         []
44         ...
45         ...
46         ...
47
48 if user_config_helper.cmd_option_exists("--help") :
49     print(USAGE)
50 else :
51     captioning = Captioning()
52     captioning.initialize()
53     speech_recognizer_data = captioning.speech_recognizer_from_user_config
54     ()
55     captioning.recognize_continuous(
56         speech_recognizer=speech_recognizer_data["speech_recognizer"],
57         format=speech_recognizer_data["audio_stream_format"], callback=
58         speech_recognizer_data["pull_input_audio_stream_callback"],
59         stream=speech_recognizer_data["pull_input_audio_stream"]
60     )
61     captioning.finish()

```

從上面的 Code 可以看出，實際的工作都由 Captioning 這個 Class 執行，然後分析這個 Class 的 Constructor 得知，程式執行縮需所有的參數都是由 user_config_helper.user_config_from_args(USEAGE) 這一個函數來的，一樣從範例給出這個函數的原始碼：

```

1 #
2 # Copyright (c) Microsoft. All rights reserved.
3 # Licensed under the MIT license. See LICENSE.md file in the project root
4 #   for full license information.
5 #
6 # 顧及篇幅，這裡把焦點放在傳回的東西上面
7 # 詳細的使用者輸入參數處理流程請參閱: https://github.com/Azure-Samples/cognitive-services-speech-sdk/blob/master/scenarios/python/console/captioning/user\_config\_helper.py
8 # file: user_config_helper.py
9
10 from datetime import timedelta
11 from enum import Enum
12 from os import linesep, environ
13 from sys import argv
14 from typing import List, Optional

```

```

15 import azure.cognitiveservices.speech as speechsdk # type: ignore
16 import helper
17
18 def user_config_from_args(usage : str) -> helper.Read_Only_Dict :
19     # start process
20     ...
21     ...
22     ...
23     # end process
24
25     return helper.Read_Only_Dict({
26         "use_compressed_audio" : cmd_option_exists("--format"),
27         "compressed_audio_format" : get_compressed_audio_format(),
28         "profanity_option" : get_profanity_option(),
29         "language" : get_language(),
30         "input_file" : get_cmd_option("--input"),
31         "output_file" : get_cmd_option("--output"),
32         "phrases" : get_phrases(),
33         "suppress_console_output" : cmd_option_exists("--quiet"),
34         "captioning_mode" : captioning_mode,
35         "remain_time" : td_remain_time,
36         "delay" : td_delay,
37         "use_sub_rip_text_caption_format" : cmd_option_exists("--srt"),
38         "max_line_length" : int_max_line_length,
39         "lines" : int_lines,
40         "stable_partial_result_threshold" : get_cmd_option("--threshold"),
41         "subscription_key" : key,
42         "region" : region,
43     })

```

從上面的 Code 可以看出，使用者輸入的參數在經過一連串的處理後，傳回一個自訂類別 `helper.Read_Only_Dict`，顧名思義，這是一個唯讀 Dict，可以保護參數在執行時不受影響。在了解上面的事情後，就可以著手改造 Code，使其能融合進專案裡面了。最後的結果如下：

```

1  # 顧及篇幅，只給出實際更動的部分。
2  # 詳細的Code請直接到captioning.py觀看
3
4  from datetime import timedelta
5  from datetime import time
6  from more_itertools import pairwise
7  from os import linesep, remove
8  from os.path import exists
9  from time import sleep
10 from typing import List, Optional
11 import wave
12 import azure.cognitiveservices.speech as speechsdk # type: ignore
13 import caption_helper
14 import helper
15 import user_config_helper
16 from user_config_helper import CaptioningMode, get_phrases,

```

```

    get_compressed_audio_format
17 from utils import AzureSpeechConfig, Logger
18 from CAPTIONINGUSAGE import CAPTIONINGUSAGE
19
20
21 class Captioning(object):
22     def __init__(self, input, output, language):
23
24         # self._user_config = user_config_helper.user_config_from_args(
25             CAPTIONINGUSAGE)
26
27         self._user_config = helper.Read_Only_Dict({
28             "use_compressed_audio": False,
29             "compressed_audio_format": get_compressed_audio_format(),
30             "profanity_option": speechsdk.ProfanityOption.Raw,
31             "language": language,
32             "input_file": input,
33             "output_file": output,
34             "phrases": get_phrases(),
35             "suppress_console_output": False,
36             "captioning_mode": CaptioningMode.OFFLINE,
37             "remain_time": timedelta(milliseconds=1000),
38             "delay": timedelta(milliseconds=1000),
39             "use_sub_rip_text_caption_format": True,
40             "max_line_length": helper.DEFAULT_MAX_LINE_LENGTH_SBCS,
41             "lines": 1,
42             "stable_partial_result_threshold": '3',
43             "subscription_key": AzureSpeechConfig.subscription_key.value,
44             "region": AzureSpeechConfig.region.value,
45         })
46         self._srt_sequence_number = 1
47         self._previous_caption: Optional[caption_helper.Caption] = None
48         self._previous_end_time: Optional[time] = None
49         self._previous_result_is_recognized = False
50         self._recognized_lines: List[str] = []
51         self._offline_results: List[speechsdk.SpeechRecognitionResult] =
52             []
53         ...
54         ...
55         ...
56
57 @Logger
58 def makeSubtitles(filePath, outputPath, lang):
59     if user_config_helper.cmd_option_exists("--help"):
60         print(CAPTIONINGUSAGE)
61     else:
62
63         captioning = Captioning(
64             input=filePath,
65             output=outputPath,

```



```

64         language=lang
65     )
66     captioning.initialize()
67     speech_recognizer_data = captioning.
        speech_recognizer_from_user_config()
68     captioning.recognize_continuous(
69         speech_recognizer=speech_recognizer_data["speech_recognizer"],
70         format=speech_recognizer_data["audio_stream_format"],
71         callback=speech_recognizer_data["
            pull_input_audio_stream_callback"],
72         stream=speech_recognizer_data["pull_input_audio_stream"]
73     )
74     captioning.finish()

```

3.3 翻譯得到的字幕檔

處理好字幕的部分之後，接下來就是翻譯了，不過翻譯時會連同時間軌的半形符號一起轉成全形符號，如果就這樣拿去給 ffmpeg 合成的話會出問題的！於是在翻譯結束後把音檔中被轉成全形的符號轉回半形就是首要任務了。這裡用 regex expression 去制定格式，當符合格式時才將全形符號轉回半形。完整的 Code 如下：

```

1  import re
2  import requests
3  import uuid
4  from utils import AzureTranslateConfig
5  from utils import Logger
6
7
8  @Logger
9  def translate(text, resource, to, outputFilePath):
10     # Add your key and endpoint
11     key = AzureTranslateConfig.subscription_key.value.decode('utf-8')
12     endpoint = AzureTranslateConfig.endpoint.value.decode('utf-8')
13
14     # location, also known as region.
15     # required if you're using a multi-service or regional (not global)
16     # resource. It can be found in the Azure portal on the Keys and
17     # Endpoint page.
18     location = AzureTranslateConfig.region.value.decode('utf-8')
19
20     path = '/translate'
21     constructed_url = endpoint + path
22
23     params = {
24         'api-version': '3.0',
25         'from': resource,
26         'to': to
27     }
28
29     headers = {

```

```

28     'Ocp-Apim-Subscription-Key': key,
29     # location required if you're using a multi-service or regional (
        not global) resource.
30     'Ocp-Apim-Subscription-Region': location,
31     'Content-type': 'application/json',
32     'X-ClientTraceId': str(uuid.uuid4())
33 }
34
35 # You can pass more than one object in body.
36 body = [{
37     'text': text
38 }]
39
40 response = requests.post(
41     constructed_url,
42     params=params,
43     headers=headers,
44     json=body
45 ).json()
46
47 with open(outputFilePath, 'w', encoding='utf-8') as f:
48     lines = response[0]['translations'][0]['text'].split('\n')
49     result = []
50     for line in lines:
51         if re.search('^\\d+:\\d+:\\d+,\\d+ --> \\d+:\\d+:\\d+,\\d+$',
52             line):
53             line = line.replace(':', ',')
54             line = line.replace(' ', ',')
55             result.append(line)
56     f.write('\\n'.join(result))

```

3.4 把字幕跟影片結合

最後，就是要把字幕跟影片結合在一起了，這一步一樣可以用 ffmpeg 完成，要在 ffmpeg 中將影片跟字幕結合，需要用到 ffmpeg 的一種機制叫做 filter，filter 簡單來說就是 ffmpeg 用來處理圖片與影音編輯的工具，舉例來說處理字幕的 filter 叫做 subtitles，有了它我們就能將套字幕這件事完美融合進專題了！

不過這個 filter 輸出的影片檔是沒有聲音的，因此我們前面抽出來的音檔又能派上用場了，只要分別定義好音檔跟經過 subtitles filter 處理的影片檔，最後融合在一起就可以了。下面是完整的程式：

```

1 import ffmpeg
2 from utils import Logger
3
4 @Logger
5 def combine_video(filepath, filename, subtitleFilePath):
6     video = (
7         ffmpeg.input(filepath)
8         .filter('subtitles', subtitleFilePath)
9     )

```

```

10
11     audio = (
12         ffmpeg.input(f'./statics/{filename}.wav')
13     )
14
15     (
16         ffmpeg.concat(video, audio, v=1, a=1)
17         .output(f'./result/{filename}_subtitled.mp4')
18         .overwrite_output()
19         .run()
20     )

```

3.5 Log 與 Decorator

為了了解程式執行的過程與每個步驟的執行時間，需要在每個函數的開頭和結尾加入以下程式：

```

1 t1 = time.time()
2 # begin program
3 ...
4 ...
5 ...
6 # end program
7 t2 = time.time()
8 print(f"cost {t2-t1}s")

```

如果要計時的函數多起來的話，那可是一項不小的功夫，於是 Decorator[4] 就派上用場了，只要寫好一個通用的計時 Decorator，將要計時的函數傳入 Decorator 中就可以幫你計時了！就像下面的 Code 所示：

```

1 import time
2 import logging
3
4
5 def Logger(func):
6     def wrap(*args, **kwargs):
7         logger = logging.getLogger()
8         logger.setLevel(logging.DEBUG)
9         formatter = logging.Formatter(
10             '[(levelname)s %(asctime)s] %(message)s',
11             datefmt='%Y%m%d %H:%M:%S'
12         )
13
14         logStream = logging.StreamHandler()
15         logStream.setLevel(logging.DEBUG)
16         logStream.setFormatter(formatter)
17         logger.addHandler(logStream)
18         logging.info(f'function {func.__name__}() start')
19         t1 = time.time()
20         result = func(*args, **kwargs)
21         t2 = time.time()

```

```

22     logging.info(f'function {func.__name__}() end')
23     logging.info(f'cost {t2-t1} s')
24     logger.removeHandler(logStream)
25     return result
26 return wrap
27
28
29 @ Logger
30 def func1():
31     print("func1 ! ")
32     time.sleep(5)
33
34
35 @ Logger
36 def func2():
37     print("func2 ! ")
38     time.sleep(10)
39
40
41 func1()
42 func2()

```

在上面的 Code 中，除了計時之外，在輸出的時候選擇了 Logging[5] 這個套件，可以更好的顯示所需要的資訊，執行結果如下：

```

1 [INFO 20221213 14:45:48] function func1() start
2 func1 !
3 [INFO 20221213 14:45:53] function func1() end
4 [INFO 20221213 14:45:53] cost 5.001032590866089 s
5 [INFO 20221213 14:45:53] function func2() start
6 func2 !
7 [INFO 20221213 14:46:03] function func2() end
8 [INFO 20221213 14:46:03] cost 10.005817413330078 s

```

對了，像 @Decorator 的寫法是 python 提供的語法糖 (Syntactic sugar)，如果不想用這個方法寫的話，可以用下面的方式使用：

```

1 Logger(func1)()
2 Logger(func2)()

```

4 後記

本次專題的不足之處在於對於歌曲的處理很糟糕，畢竟每首歌的結構與狀況都不一樣，實在很難有一個比較統一的標準，難度自然就往上升了，而對於一般的影片效果相對還是不錯的。

藉由這次專題，對於 Azure 可以說更加熟悉了，相對於為了考過 AI900 背的題目，果然還是實作更有趣呢。對了，結果放在 result 資料夾中。

參考資料

- [1] Magiclen. FFmpeg：免費開源、功能強大的影音處理框架. URL: <https://magiclen.org/ffmpeg/>.
- [2] kkroening. Python bindings for FFmpeg - with complex filtering support. URL: <https://pypi.org/project/ffmpeg-python/>.
- [3] Microsoft. 快速入門：使用語音轉換文字建立字幕. URL: <https://learn.microsoft.com/zh-tw/azure/cognitive-services/speech-service/captioning-quickstart?tabs=windows%2Cterminal&pivots=programming-language-python>.
- [4] codemee. Python 裝飾器 (Decorator). URL: <https://dev.to/codemee/python-decorator-3bni>.
- [5] shengyu7697. Python logging 日誌用法與範例. URL: <https://shengyu7697.github.io/python-logging/>.