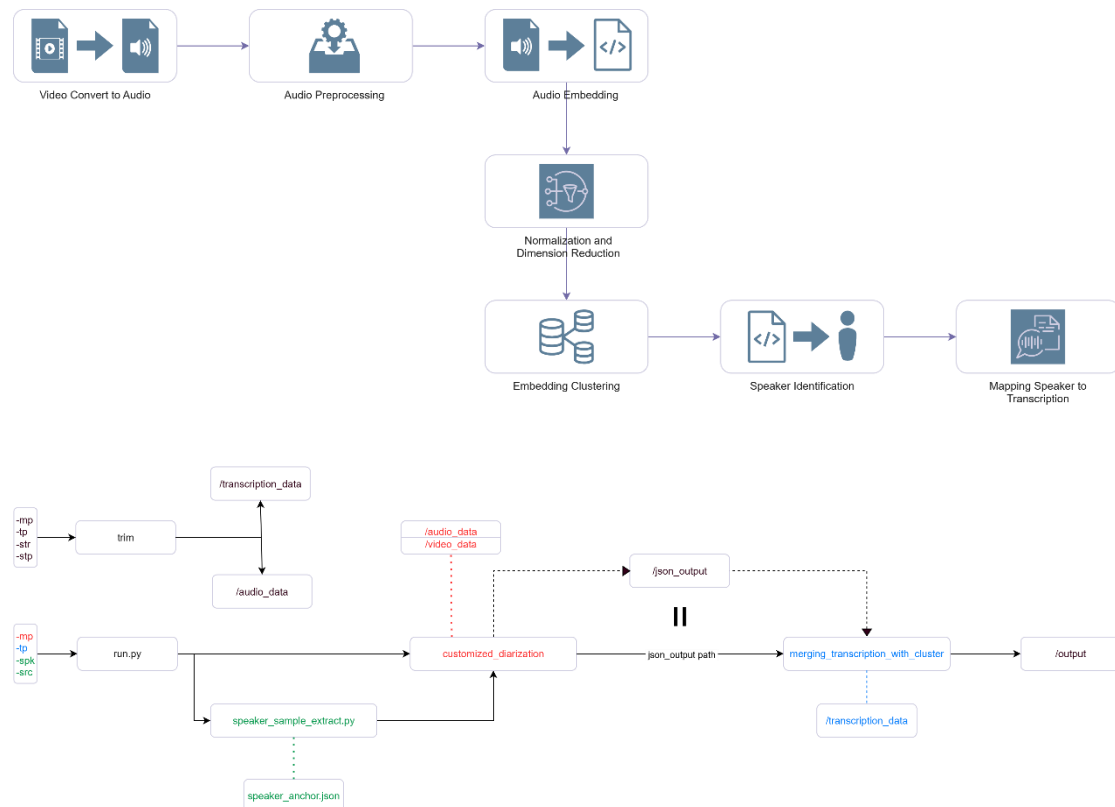


# Speaker Diarization

## Environment

- Operation System: Microsoft Windows 11 Pro
- CPU: 13th Gen Intel(R) Core(TM) i5-1335U
- GPU: NVIDIA GeForce RTX 2050
- Python: 3.11.9 (Using pyenv + venv)

## Workflow



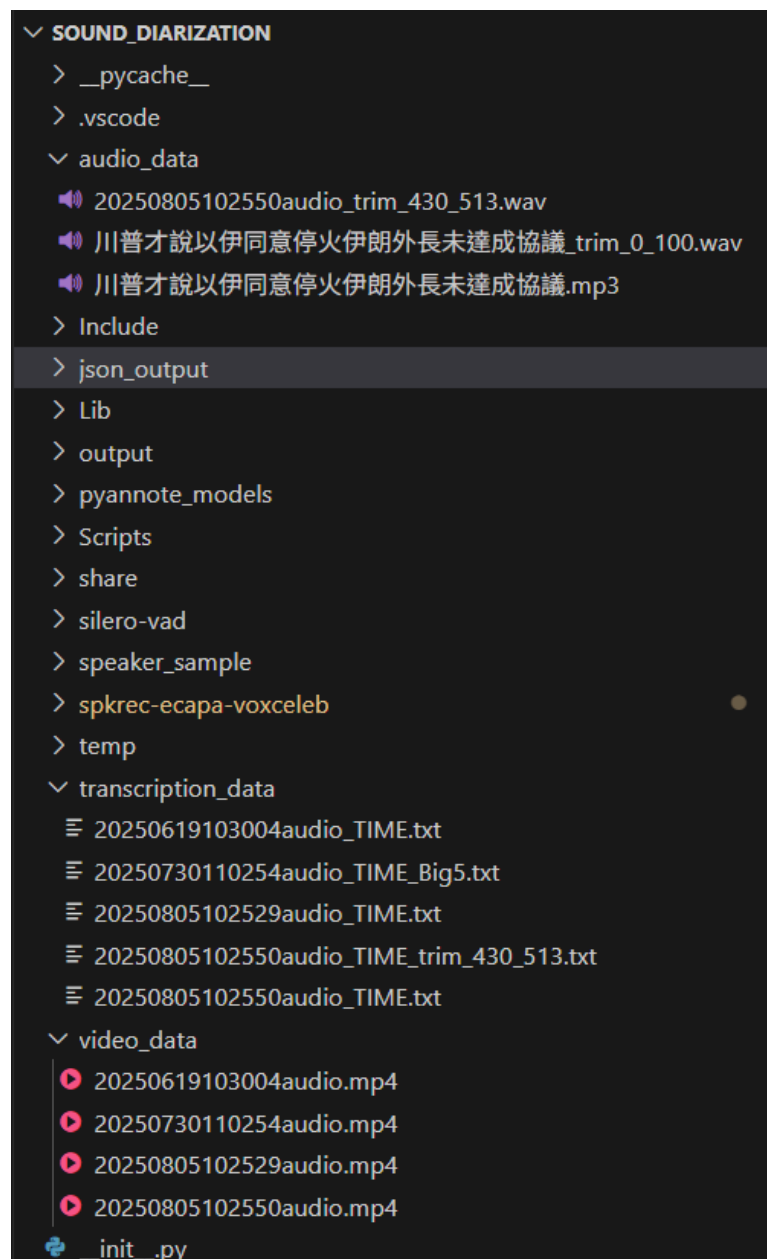
# File Hierarchy

```
├── /sound_diarization
│   ├── pyannotate_diarization_config.json
│   ├── customized_diarization_config.json
│   ├── speaker_anchor.json
│   ├── __init__.py
│   ├── run.py
│   ├── pyannotate_diarization_offline.py
│   ├── customized_diarization.py
│   ├── speaker_sample_extract.py
│   ├── merging_transcription_with_cluster.py
│   ├── utils.py
│   ├── requirements.txt
│   ├── /pyannotate_models
│   │   ├── ...
│   │   └── pyannotate_diarization_config.yaml
│   ├── /silero-vad
│   │   ├── ...
│   │   └── hubconf.py
│   ├── /spkrec-ecapa-voxceleb
│   │   ├── ...
│   │   └── hyperparams.yaml
│   ├── /audio_data
│   │   ├── ...
│   │   └── test.mp3
│   ├── /video_data
│   │   ├── ...
│   │   └── test.mp4
│   ├── /transcription_data
│   │   ├── ...
│   │   └── test.txt
│   ├── /speaker_sample
│   │   ├── ...
│   │   ├── speaker1.npy
│   │   └── speaker1.wav
│   ├── /temp
│   ├── /json_output
│   └── /output
```

# Data

This project has 3 types of input data:

- Audio Data: The audio file can be any audio format. All audio data will be stored in the folder `audio\_data`.
- Video Data: The video file can be any video format. All video data will be stored in the folder `video\_data`.
- Transcription Data: The transcription file should be a .txt file with a strict row-by-row format; `[(start:.2f)s -> (stop:.2f)s] sentences`. All transcription data will be stored in the folder `transcription\_data`.

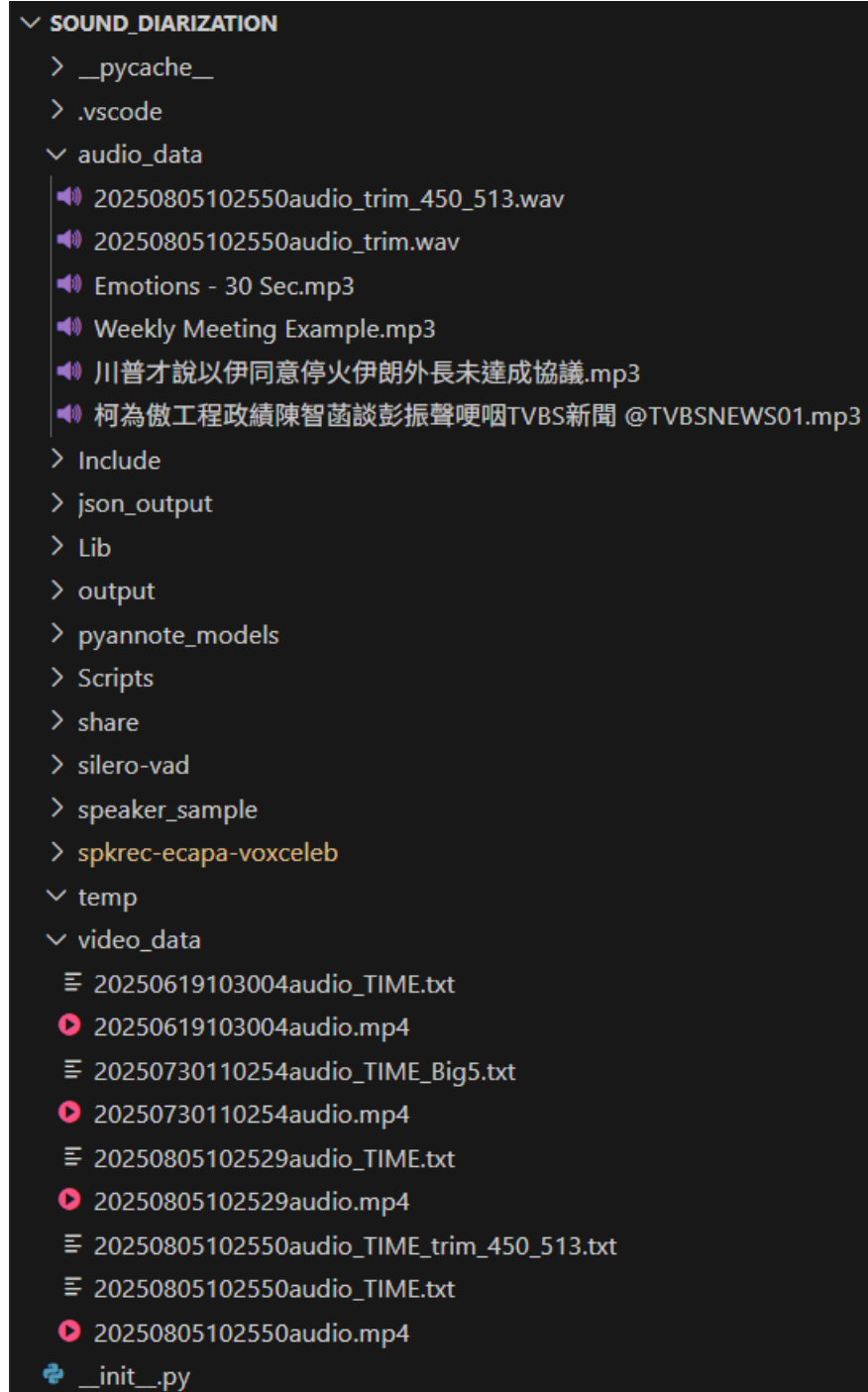


# Introduction

This package includes 2 diarization programs:

- `pyannote_diarization_offline`:  
This program uses the `pyannote/speaker-diarization-3.1` model to perform speaker diarization, this program serves as the benchmark for our `customized_diarization` program.
- `customized_diarization`:  
This program performs speaker diarization by implementing the whole workflow from scratch. It supports detailed model tuning and customized parameters adjustment at each stage of the workflow.

Each program supports audio or video as the input file. The audio files should be placed in folder `audio_data`, and the video files should be placed in folder `video_data`.



## pyannote\_diarization\_offline

### pyannote\_diarization\_config.json

The JSON file `pyannote\_diarization\_config.json` controls what kind of preprocessing techniques will apply to the input file.

```

{} pyannote_diarization_config.json > ...
1  {
2    "preprocessing":{
3      "vad": true,
4      "vad_keep_length": false,
5      "pre_emphasis": true
6    }
7  }

```

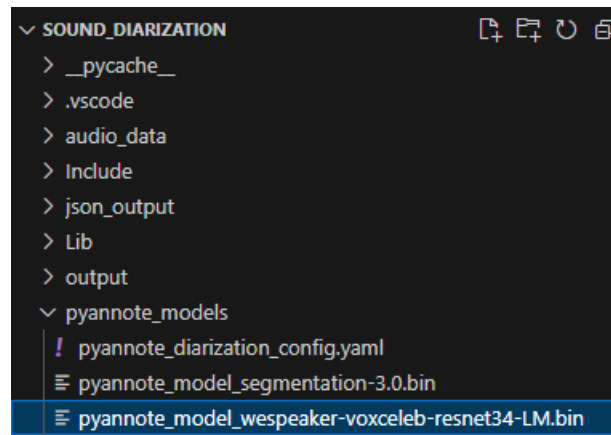
Parameter	Data Type	Description
vad	Bool	If True, voice activity detection will apply. Default is True.
vad_keep_length	Bool	If True, vad will remove the non-speech portions and concatenate the speech portions. Otherwise, the non-speech portions will be replaced with silence. Default is False.
pre_emphasis	Bool	If True, Pre-emphasis will apply. Default is True.

## / pyannote\_models

The folder `pyannote\_models` contains the necessary files for deploying the speaker diarization model `pyannote/speaker-diarization-3.1` on-premises.

To know more details about the implementation, see

[https://github.com/pyannote/pyannote-audio/blob/main/tutorials/community/offline\\_usage\\_speaker\\_diarization.ipynb](https://github.com/pyannote/pyannote-audio/blob/main/tutorials/community/offline_usage_speaker_diarization.ipynb)



## customized\_diarization

### run

The program `run.py` is the main program of the customized\_diarization workflow.

This program accepts 4 parameters:

```
parser.add_argument("--mp", "--media_path", help="input media path", dest="media_path", type=str, required=True)
parser.add_argument("--tp", "--transcription_path", help="transcription path", dest="transcription_path", type=str, required=False)
parser.add_argument("--spk", "--speaker", help="The string assigned speakers join with '|', e.g. 'Jeremy|Barry'", dest="speaker", type=str, required=False)
parser.add_argument("--src", "--source", help="The string assigned source join with '|', e.g. './file1.mp4|./file2.mp4'", dest="source", type=str, required=False)
```

## customized\_diarization\_config.json

The JSON file `customized\_diarization\_config.json` controls the detailed parameters at each stage applying to the input file.

```

1  {} customized_diarization_config.json > {} dimension_reduction > criteria
2  {
3      "init": {
4          "embed": "ecapa",
5          "reduction": "PCA",
6          "normalization": true,
7          "cluster": "HDBSCAN",
8          "speakers": 5,
9          "threshold": 0.65
10     },
11     "segmentation": {
12         "window": 1.6,
13         "stride": 0.08,
14         "batch": 1
15     },
16     "preprocessing": {
17         "vad": true,
18         "vad_keep_length": false,
19         "pre_emphasis": true
20     },
21     "dimension_reduction": {
22         "ratio": 0.85,
23         "n_components": 25,
24         "criteria": "ratio",
25         "random_state": 42,
26         "kpca": {
27             "kernel": "rbf",
28             "gamma": 1
29         }
30     },
31     "clustering": {
32         "hdbscan": {
33             "min_cluster_size": 32,
34             "min_samples": 32,
35             "metric": "cosine",
36             "alpha": 0.5
37         },
38         "affinitypropagation": {
39             "random_state": 5,
40             "convergence_iter": 50,
41             "max_iter": 1000,
42             "damping": 0.5,
43             "distance": "euclidean"
44         },
45         "spectralclustering": {
46             "random_state": 5,
47             "affinity": "rbf",
48             "gamma": 1.0
49         }
50     }
51 }

```

Parent	Parameter	DataType	Description
init	embed	Str	The embedding method. Default is “ecapa”.
init	reduction	Str	The dimension reduction method. It can be “pca” or “kpca”, with a default value of “pca”.
init	normalization	Bool	If True, the L2 normalization will apply to the audio embeddings. Default is True.
init	cluster	Str	The clustering method. It can



			be “HDBSCAN”, “AP” (affinity propagation) or “SC” (spectral clustering), with a default value of “HDBSCAN”.
init	speakers	Int	The number of speakers present in the input file, only needs to be specified if “SC” is chosen as the clustering method.
init	threshold	Float	The threshold to consider that a segment embedding belongs to a certain speaker, the value range is between 0 and 1, with a default value of 0.65.
segmentation	window	Float	The window size of a segment in seconds. Default is 1.6.
segmentation	stride	Float	The stride size of a segment in seconds. Default is 0.08.
segmentation	batch	Int	The batch size of a segment. Default is 1.
preprocessing	vad	Bool	If True, voice activity detection will apply. Default is True.
preprocessing	vad_keep_length	Bool	If True, vad will remove the non-speech portions and concatenate the speech portions. Otherwise, the non-speech portions will be replaced with silence. Default is False.
preprocessing	pre_emphasis	Bool	If True, Pre-emphasis will apply. Default is True.
dimension_reduction	ratio	Float	The expected explained variance ratio used to determine how many principal components to retain. Default is 0.85.

dimension_reduction	n_components	Int	The number of components that will yield after dimension reduction. Default is 25.
dimension_reduction	criteria	Str	Specifies whether to apply `n_components` or `ratio` during dimensionality reduction. Default is `n_components`.
dimension_reduction	random_state	Int	The random state for the dimension reduction methods.
dimension_reduction	kpca	Dict	The specific parameters for `KPCA` method including `kernel` and `gamma`. See <a href="https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.KernelPCA.html">https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.KernelPCA.html</a> to get more info.
clustering	hdbscan	Dict	The specific parameters for `HDBSCAN`. See <a href="https://scikit-learn.org/stable/modules/generated/sklearn.cluster.HDBSCAN.html">https://scikit-learn.org/stable/modules/generated/sklearn.cluster.HDBSCAN.html</a> to get more info.
clustering	affinitypropagation	Dict	The specific parameters for `affinitypropagation`. The parameter `distance` isn't from the original scikit-learn docs, this parameter determines the kind of distance applied to calculating the similarity of datapoints For other parameters, see <a href="https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AffinityPropagation.html">https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AffinityPropagation.html</a>

			<a href="https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AffinityPropagation.html">learn.org/stable/modules/generated/sklearn.cluster.AffinityPropagation.html</a> to get more info.
clustering	spectralclustering	Dict	The specific parameters for `spectralclustering`. See <a href="https://scikit-learn.org/stable/modules/generated/sklearn.cluster.SpectralClustering.html">https://scikit-learn.org/stable/modules/generated/sklearn.cluster.SpectralClustering.html</a> to get more info.

### speaker\_anchor.json

The JSON file `speaker\_anchor.json` stores a list speaker info, where each item defines a specific speaker's audio clip in a video/audio file.

Parameter	Data Type	Description
name	Str	The name of the speaker.
source	Str	The file path of the media file (audio or video) from which the segment is taken.
start	Int	The start time of the speaker's segment in seconds.
stop	Int	The stop time of the speaker's segment in seconds.

```

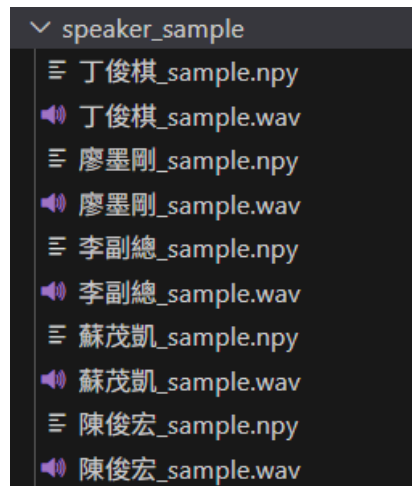
{} speaker_anchor.json > [ ] speaker
2      "speaker": [
51      {
52          "name": "陳俊宏",
53          "source": "./video_data/20250805102550audio.mp4",
54          "start": 300,
55          "stop": 330
56      },
57      {
58          "name": "李副總",
59          "source": "./video_data/20250805102550audio.mp4",
60          "start": 408,
61          "stop": 430
62      },
63      {
64          "name": "蔡昀儒",
65          "source": "./video_data/20250805102529audio.mp4",
66          "start": 183,
67          "stop": 217
68      },
69      {
70          "name": "蘇茂凱",
71          "source": "./video_data/20250805102529audio.mp4",
72          "start": 293,
73          "stop": 323
74      }
75  ]
76  }

```

## /speaker\_sample

The folder `speaker\_sample` stores speaker sample embeddings and clips generating by `speaker\_sample\_extract.py`.

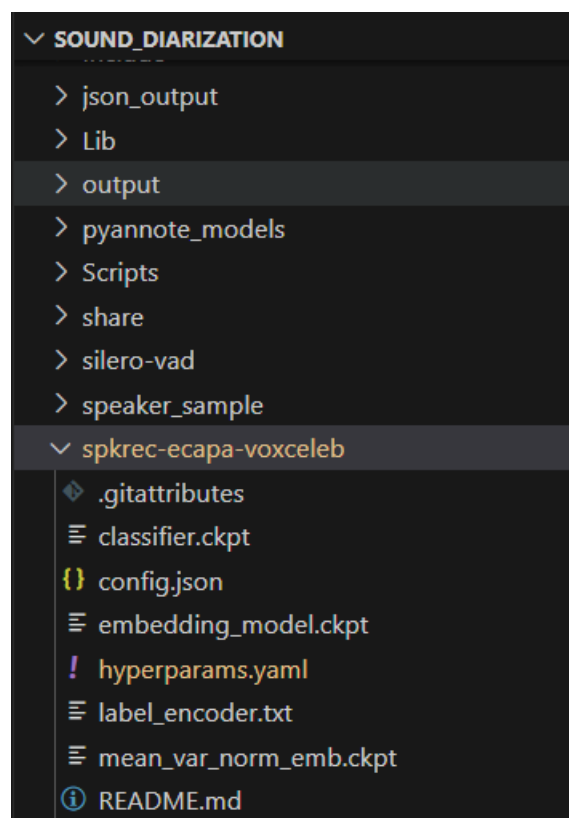
- Speaker sample embeddings are saved as .npy file.
- Speaker sample clips are saved as .wav file.



## / spkrec-ecapa-voxceleb

The folder `spkrec-ecapa-voxceleb` contains the ECAPA-TDNN speaker embedding model repository. We download the entire repository and place it under the root directory of the environment to enable the on-premises speaker diarization embedding workflow.

The link of the repository: [speechbrain/spkrec-ecapa-voxceleb · Hugging Face](https://huggingface.co/speechbrain/spkrec-ecapa-voxceleb)



## /temp

The folder `temp` stores intermediate files, which will be removed once the

workflow complete.

## speaker\_sample\_extract

The program `speaker\_sample\_extract.py` generates the speaker sample embeddings and clips, these outputs are saved in the folder `speaker\_sample`.

## merging\_transcription\_with\_cluster

The program `merging\_transcription\_with\_cluster.py` merges the transcription with the corresponding speakers/clusters line by line.

## utils

The Python file `utils.py` contains all the general-purpose functions used throughout this project.

## trim

In some cases, the input media file may be too large to process, potentially causing the computer to shut down due to memory exhaustion. To address this, a simple function is provided to trim or extract a small clip of the media file along with the corresponding transcription.

```
$ python trim.py -mp "./video_data/20250805102550audio.mp4" -tp "video_data\20250805102550audio_TIME.txt" -str 430 -stp 513
trimming ./video_data/20250805102550audio.mp4 from 430s to 513s...
MoviePy - Writing audio in temp\20250805102550audio.wav
MoviePy - Done.

trimming video_data\20250805102550audio_TIME.txt from 430s to 513s...
[1.64s -> 5.36s] 應該是說廠商他在他場內測試是OK的資料上傳
[5.36s -> 7.16s] 還有建一個
[7.16s -> 10.88s] 就是他們自己的虛擬的資料庫去模擬上傳動作
[11.60s -> 14.52s] 問題是那個就是一個虛擬的環境
[14.52s -> 16.52s] 這裡就是戰場就是在這裡
[16.52s -> 18.72s] 他在家裡面去搞一個
```

The `trim` program accepts 4 parameters:

1. `media_path`: The path to the media file, which can be either a video or an audio file.  
`transcription_path`: The path to the transcription file, which must be a .txt file with a strict row-by-row format; `[(start:.2f)s -> (stop:.2f)s] sentences`.

**regular expression:** `r"[(\d+\.\d{2})s -> (\d+\.\d{2})s] \w+\n"`

2.

```
10 [76.68s -> 79.60s] 以及偵測履歷的部分會去做展開說明
11 [79.60s -> 82.24s] 那還有兩次數位化目前整個專案的推動進度
12 [82.24s -> 85.04s] 那另外在針對AI視覺辨識的話
13 [85.04s -> 90.48s] 我們這邊本次有針對於這一兩個月的新東西做推進的說明
14 [90.48s -> 94.64s] 那等一下最後還有針對於那個公司數位化發展部分
```

3. start: The start time (in seconds) of the trimmed segment.

4. stop: The stop time (in seconds) of the trimmed segment.

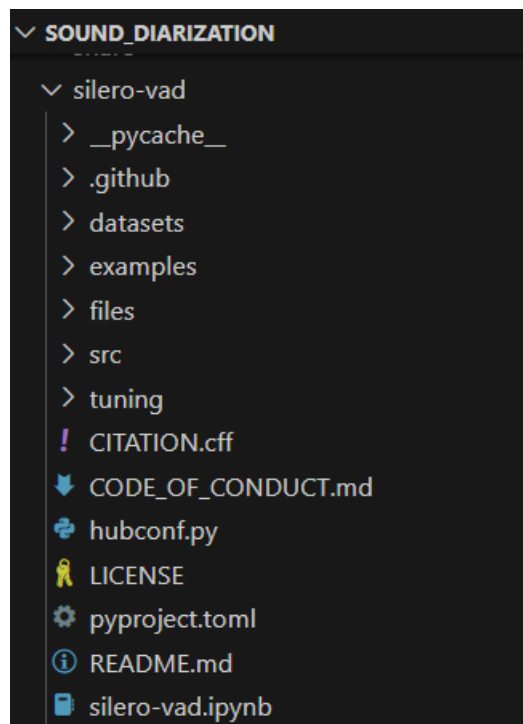
Regardless of whether the input is an audio or video file, the clips will be exported as a .wav file in the `audio\_data` folder, and the corresponding transcription will be saved as a .txt file in the `video\_data` folder.

## / silero-vad

The folder `silero-vad` contains the Voice Activity Detection (VAD) repository. We download the entire repository and place it under the root directory of the environment to enable on-premises speaker diarization VAD, which is utilized by both `pyannote_diarization_offline` and `customized_diarization`.

The link to the repository and implementation:

<https://github.com/snakers4/silero-vad/blob/master/silero-vad.ipynb>

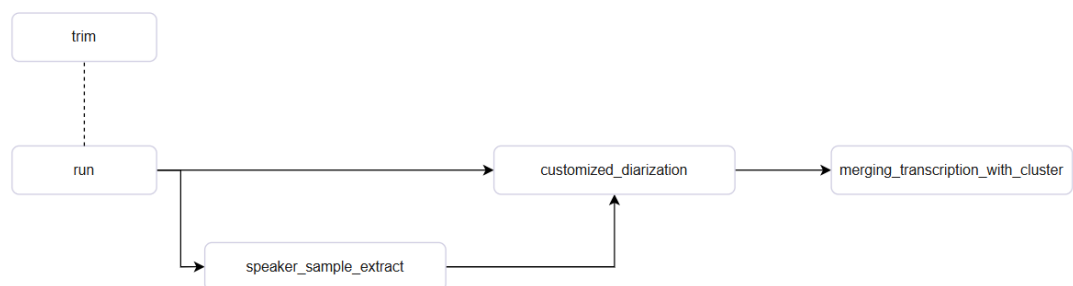


## Usage

### pyannote\_diarization\_offline

```
$ python pyannote_diarization_offline.py -p "audio_data\\川普才說以伊同意停火伊朗外長未達成協議.mp3"
Loading pyannote pipeline from pyannote_models\pyannote_diarization_config.yaml...
Changing working directory to C:\Users\680451\Documents\Work Space\sound_diarization
Changing working directory back to C:\Users\680451\Documents\Work Space\sound_diarization
start=0.03s stop=30.96s speaker_SPEAKER_01
start=31.05s stop=62.76s speaker_SPEAKER_02
start=62.76s stop=71.92s speaker_SPEAKER_03
start=72.14s stop=95.00s speaker_SPEAKER_02
start=95.07s stop=105.90s speaker_SPEAKER_00
start=106.04s stop=124.38s speaker_SPEAKER_02
start=124.62s stop=132.96s speaker_SPEAKER_04
start=133.14s stop=143.92s speaker_SPEAKER_02
```

### customized\_diarization





```

$ python run.py -mp "audio_data/20250805102550audio_trim_430_513.wav" -tp "transcription_data/20250805102550audio_TIME_trim_430_513.txt" -spk "陳俊宏|李副總|蘇茂凱"
[{'name': '蘇茂凱', 'source': './video_data/20250805102529audio.mp4', 'start': 293, 'stop': 323}, {'name': '陳俊宏', 'source': './video_data/20250805102550audio.mp4', 'start': 300, 'stop': 330}, {'name': '李副總', 'source': './video_data/20250805102550audio.mp4', 'start': 408, 'stop': 430}]
Speaker Embedding 蘇茂凱_sample.npy has already existed.
Speaker Embedding 陳俊宏_sample.npy has already existed.
Speaker Embedding 李副總_sample.npy has already existed.
temp\20250805102550audio_trim_430_513.wav
100% | 823/823 [01:53:00:00, 7.24it/s]
embedding array shape: (823, 192)
Normalization: True
dimension reduction method: PCA
general_params: {'ratio': 0.85, 'n_components': 25, 'criteria': 'n_components', 'random_state': 42}
clustering method: HDBSCAN
HDBSCAN:
{'min_cluster_size': 32, 'min_samples': 32, 'metric': 'cosine', 'alpha': 0.5}
before: 0, after: 10
before: 10, after: 6
mapping clusters to transcription...
(sound diarization)

```

## Output

In the clustering stage, the model will assign a cluster to each segment. And the subsequent down-stream processing flow will combine adjacent segments having the same cluster and apply other processes to smooth the output result.

There are 3 types of cluster labels:

1. Positive Number:

The original cluster index output directly from the clustering model.

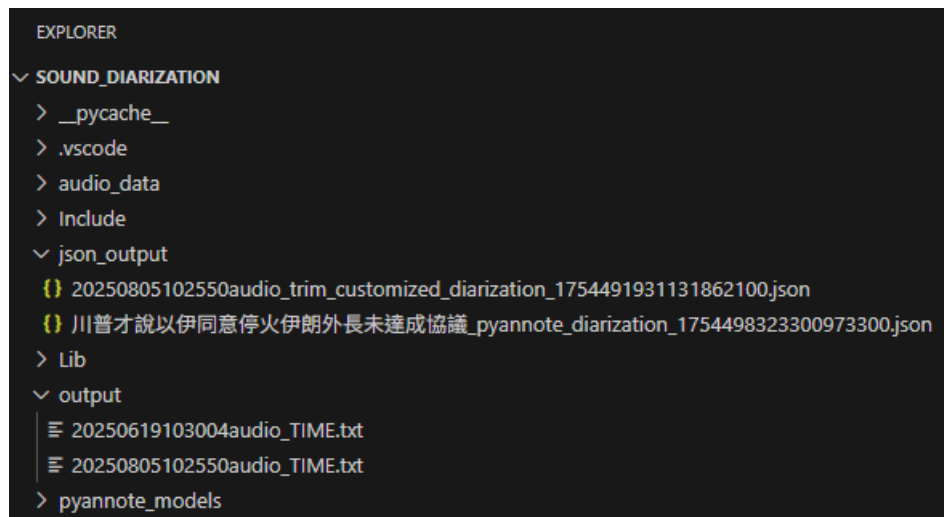
2. Name:

When a transcription file is provided and the cosine similarity score between the speaker embedding and the segment embedding exceeds the threshold specified in `customized\_diarization\_config.json`, the original cluster index will be replaced by speaker label with the highest similarity score.

3. -1:

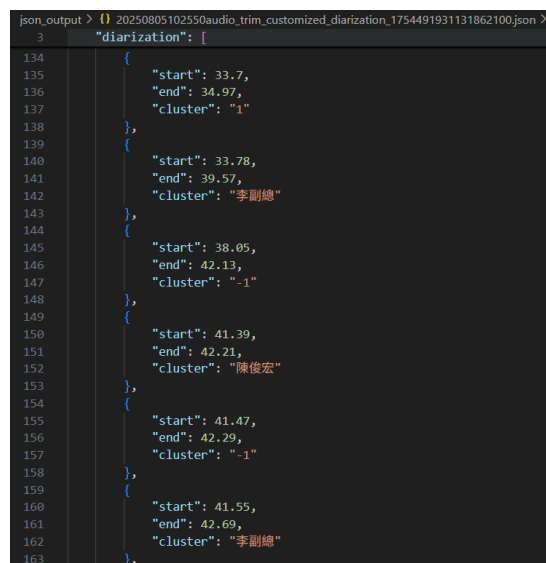
Noise, or unrecognized voice.

There are 2 different output formats in this package.



- JSON file:

The JSON format includes clusters along with their timelines, indicating start and end times in seconds. These files are stored in the folder `json\_output`.



- Txt file

The txt file resembles the original transcription input with the annotation for corresponding clusters or speakers in each row. These files are stored in the folder `output`.

```
1 [0.32s -> 1.48s] [李副總] 一切都很正常
2 [1.48s -> 3.76s] [李副總] 當然不會長出病毒來
3 [3.76s -> 6.20s] [李副總] 這裡就是什麼都有對不對
4 [6.20s -> 9.00s] [李副總] 當然病毒就會在這裡跑出來
5 [9.00s -> 12.68s] [李副總] 這樣子他去在家裡搞一個模擬的環境
6 [12.68s -> 16.88s] [-1] 跟實際現場打仗的狀況是一樣的嗎
7 [16.88s -> 20.28s] [-1] 好我會要求廠商就是看這禮拜最晚下禮拜
8 [20.28s -> 22.56s] [1] 看能不能入場拿去做一個實際的
9 [22.56s -> 25.96s] [-1] 因為就是他直接在場內做一個實際測試
10 [25.96s -> 27.68s] [李副總] 剛剛前面要講的問題是說
11 [27.68s -> 30.00s] [-1] 對啦現在就是我們想一個廠商對應一個
12 [30.00s -> 31.48s] [李副總] 想兩個廠商對應兩個
13 [31.48s -> 33.96s] [李副總] 廠商自己這樣好像也沒什麼招數可以拿出來
```

Since the `pyannote\_diarization\_offline` program just serves as a performance benchmark, it only yields the JSON format file and does not support the speaker identification process. That is, that it doesn't include speaker name labels.