# Speech And Non-speech Recognizer

## Executive Summary

This report provides a comprehensive overview of the **Audio Analysis Pipeline**, a modular Python-based system designed to process an input audio file and generate a structured analysis report. The pipeline separates audio into speech and non-speech components, transcribes spoken content into text, identifies background environmental sounds, and compiles everything into a unified final report.

The system is orchestrated by `main.py`, which sequentially executes three core scripts: `Seperator.py` (for source separation and denoising), `speech_analyser.py` (for speech-to-text transcription), and `nonspeech_analyser.py` (for non-speech sound classification). It leverages state-of-the-art open-source models for audio processing, including **Demucs** for separation, **Whisper** for transcription, and **YAMNet** for sound classification. The pipeline is configurable via file paths and parameters in each script, ensuring flexibility for different input audio files.
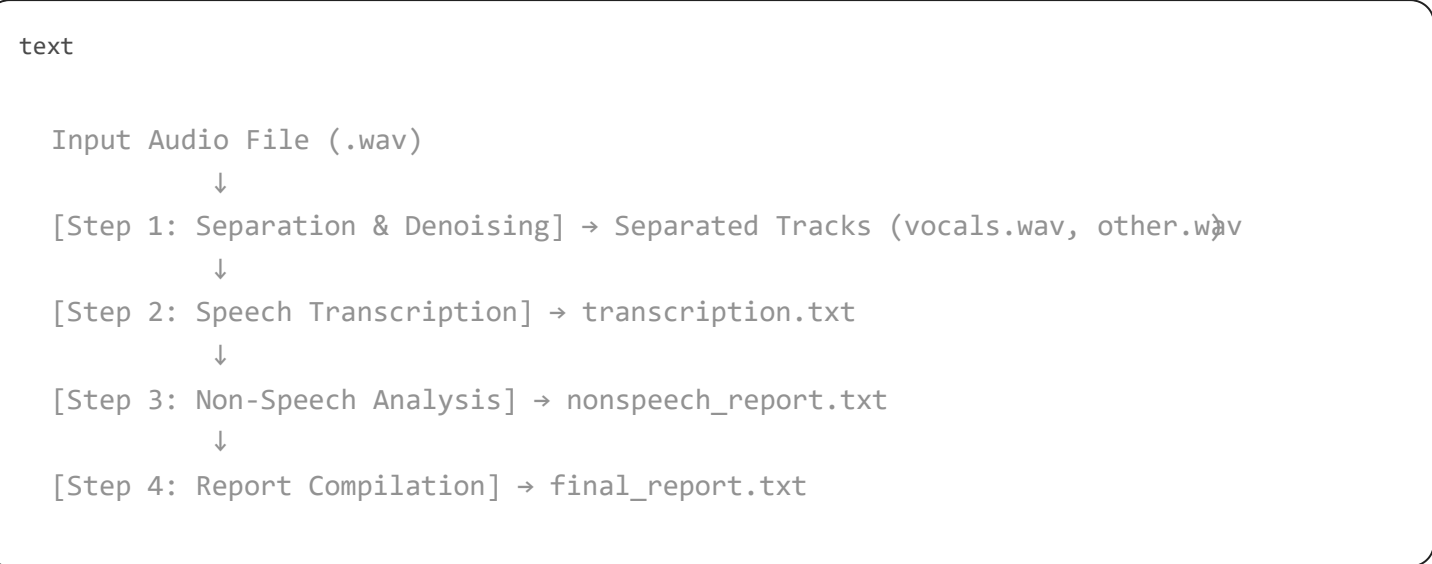
Key benefits:

- **Modularity**: Each step is isolated, allowing easy debugging and extension.
- **Efficiency**: Models are loaded once where possible, and processing runs on CPU (with potential for GPU acceleration).
- **Output**: A clean, human-readable final report in TXT format, including transcribed speech and detected non-speech sounds.

The system assumes input audio in WAV format (or convertible) at 16kHz sample rate. Processing time varies by audio length (e.g., 1-5 minutes for a 30-second clip on standard hardware).

## System Architecture and Workflow

The pipeline follows a linear, sequential workflow triggered by `main.py`. Below is a high-level diagram:

```text

  Input Audio File (.wav)
           ↓
  [Step 1: Separation & Denoising] → Separated Tracks (vocals.wav, other.wav)
           ↓
  [Step 2: Speech Transcription] → transcription.txt
           ↓
  [Step 3: Non-Speech Analysis] → nonspeech_report.txt
           ↓
  [Step 4: Report Compilation] → final_report.txt
```

Each step is detailed below, including the models used, key algorithms, inputs/outputs, and code highlights.

## Step 1: Audio Source Separation and Denoising ( `Seperator.py` )

**Purpose**: Decompose the input audio into isolated stems (e.g., vocals/speech vs. background/instruments/noise). This enables targeted analysis of speech and non-speech elements. Post-separation, noise is specifically reduced from the vocal track to improve transcription accuracy.

**Models and Libraries Used**:

- **Primary Model**: **Hybrid Transformer Demucs (htdemucs)** – A state-of-the-art neural audio source separation model from the Demucs library (by Facebook AI). It uses a hybrid architecture combining convolutional and transformer layers to separate up to 4 sources (vocals, drums, bass, other). Trained on MUSDB18 dataset for music separation but generalizes well to speech + background scenarios.
  - Why htdemucs? It outperforms older models like Spleeter in separation quality (SDR > 8dB on benchmarks) and handles overlapping signals robustly.
  - Loaded via: `get_model("htdemucs")` from `demucs.pretrained` .
- **Denoising Library**: **noisereduce** – A spectral gating-based noise reduction algorithm. It estimates noise profiles from the signal and subtracts them using Wiener filtering principles.
  - Parameters: `prop_decrease=0.9` (aggressive noise reduction), `stationary=False` (handles nonstationary noise like echoes).
- **Supporting Libraries**: `torchaudio` (for I/O and tensor handling), `torch` (PyTorch backend for Demucs), `numpy` (array operations).

**Key Algorithms**:

1. Load input audio at 16kHz (resampled if needed).
2. Apply Demucs separation: Forward pass through the model to predict source waveforms. Output: Tensor of shape `[num_sources, channels, length]` .
3. Save separated tracks (e.g., `vocals.wav` , `drums.wav` , `bass.wav` , `other.wav` ).
4. **Denoising (Vocals Only)**: Convert vocal tensor to NumPy, apply spectral noise reduction, reconvert to tensor, and overwrite the file.

**Inputs**:

- `INPUT_AUDIO` : Path to raw audio file (e.g., `D:\programs\python\speech nonspeech\4.wav` ).

**Outputs**:

- Folder: `separated/` containing WAV files for each source (e.g., `vocals.wav` for speech, `other.wav` for background).
- Console logs for progress (e.g., "Sources: ['vocals', 'drums', 'bass', 'other']").

**Code Highlights**:

```Python
# Model Loading
model = get_model("htdemucs")  # htdemucs model
model.cpu().eval() # CPU inference
```

```
    # Separation
    sources = apply.apply_model(model, wav,device="cpu")[0]  # Demucs forward pass

    # Denoising (noisereduce)
    reduced = nr.reduce_noise(y=audio_np, sr=sr, prop_decrease=0.9, stationary=False)
```

**Potential Limitations**: CPU-only by default (slow for long files >5min); assumes 4-source separation (may need customization for pure speech podcasts).

**Execution Trigger**: Called first in `main.py` via `subprocess.run()` .

## Step 2: Speech-to-Text Transcription ( `speech_analyser.py` )

**Purpose**: Transcribe the isolated vocal track into readable text, supporting multilingual audio (focus on English and Hindi). Language detection ensures model selection for accuracy.

**Models and Libraries Used**:

- **Primary Model**: **OpenAI Whisper (large-v3 variant)** – A transformer-based encoder-decoder model for automatic speech recognition (ASR). Trained on 680k hours of multilingual data, achieving <5% WER on benchmarks like LibriSpeech.
  - Variants:
    - `large` for Hindi ( `model_hi` ): Fine-tuned for Indic languages, handles code-switching.
    - `large` for English ( `model_en` ): Optimized for clean speech post-denoising.
  - Why Whisper? End-to-end (no alignment needed), robust to accents/noise, and supports 99+ languages.
  - Loaded via: `whisper.load_model("large")` .
- **Detection Model**: **Whisper base** – Lightweight variant ( `lang_model` ) for initial language detection via beam search on the first 30 seconds.
- **Supporting Libraries**: `whisper` (OpenAI's library), `os` (file handling), `datetime` (unused in provided code but imported).

**Key Algorithms**:

1. **Language Detection**: Transcribe a short clip with base model; extract `detected_lang` (e.g., "en", "hi").
2. **Transcription**: Route to appropriate large model based on language (default to English). Use beam search for decoding.
3. **Output Saving**: Clean text (no timestamps) saved as UTF-8 TXT.

**Inputs**:

- Audio path: `separated/vocals.wav` (post-separation and denoising).

**Outputs**:

- File: `Speech text output/transcription.txt` containing raw transcribed text.
- Console: Detected language and final text preview.

**Code Highlights:**

```python
Python


  # Model Loading (global for efficiency)
  model_hi = whisper.load_model("large")  # Hindi
  model_en = whisper.load_model("large")  # English
  lang_model = whisper.load_model("base")  # Detection

  # Transcription Logic
  detect = lang_model.transcribe(path)
  detected_lang = detect["language"]
  if detected_lang in ["hi", "ur", "hr"]:
      result = model_hi.transcribe(path, language="hi")
  else:
      result = model_en.transcribe(path, language="en")
  text_output = result["text"]
```

**Potential Limitations**: Assumes clean vocals (relies on prior denoising); Hindi/Urdu/Hr detection may falter on dialects; no speaker diarization.

**Execution Trigger**: Called after separation in `main.py` .

## Step 3: Non-Speech Sound Classification ( `non-speech_analyser.py` )

**Purpose**: Analyze the background track to classify environmental sounds, excluding speech-related labels. This identifies noise, music, or ambient elements for context in the report.

**Models and Libraries Used**:

- **Primary Model**: **YAMNet** – Google's pre-trained audio event classification model from TensorFlow Hub. A MobileNetV1-based CNN trained on AudioSet (2M+ clips, 521 classes). Outputs probabilities for 521 sound events at 10-second clips.

  - Loaded via: `hub.load("https://tfhub.dev/google/yamnet/1")` .

  - Class Map: 521 labels (e.g., "Dog", "Music", "Traffic noise") from AudioSet CSV.

  - Why YAMNet? Lightweight (1MB), real-time capable, and excels at weakly supervised classification (no bounding boxes needed).

- **Supporting Libraries**: `tensorflow` / `tensorflow_hub` (model inference), `librosa` (audio loading/resampling to 16kHz), `numpy` / `pandas` (score aggregation and CSV handling).

**Key Algorithms**:

1. Load background audio ( `other.wav` ) and extract waveform.

2. Inference: Run through YAMNet to get per-frame scores (shape: `[frames, 521]` ).

3. Aggregation: Mean pooling over time for average class scores.

4. Filtering: Sort top scores, exclude speech terms (e.g., "Speech", "Conversation"), limit to top 10.

5. Report Generation: TXT with label-score pairs.

**Inputs**:

- Audio path: `separated/other.wav` .

**Outputs**:

- File: `Speech text output/nonspeech_report.txt` (e.g., "Dog: 0.856\nMusic: 0.723").
- Console: Report preview.

**Code Highlights**:

```Python
# Model Loading
yamnet = hub.load("https://tfhub.dev/google/yamnet/1")
class_names = pd.read_csv(class_map_path)['display_name'].tolist()  # 521 classes

# Inference
scores, _, _ = yamnet(waveform)  # [frames, 521]
avg_scores = np.mean(scores.numpy(), axis=0)

# Top Non-Speech
speech_terms = ["Speech", "Conversation", "Narration", "Babbling"]
top_indices = np.argsort(avg_scores)[::-1]
# Filter and collect top 10
```

**Potential Limitations**: Coarse-grained (event-level, not timestamped); may misclassify music as "nonspeech" if overlapping.

**Execution Trigger**: Called last before compilation in `main.py` .

## Step 4: Report Compilation and Orchestration ( `main.py` )

**Purpose**: Coordinate the pipeline, extract cleaned outputs from prior steps, and generate a unified report. Handles errors and provides user-friendly logging.

**Models and Libraries Used**:

- None (pure orchestration). Relies on `subprocess` for script execution, `re` for text cleaning.

**Key Algorithms**:

1. Sequential Execution: Run scripts via `subprocess.run([sys.executable, script_path])` .
2. Extraction:
   - Speech: Read `transcription.txt` , clean whitespace.
   - Non-Speech: Parse `nonspeech_report.txt` , skip headers, extract label:score lines.
3. Compilation: Template-based TXT report with sections for speech and non-speech.
4. Error Handling: Check file existence; exit on failures.

**Inputs**:

- Base directory: `D:\programs\python\speech nonspeech` .

**Outputs**:

- Folder: `Final report/final_report.txt` (formatted report).
- Console: Progress logs and full report preview.

**Code Highlights**:

```python
Python

  # Orchestration
  run_script("Seperator.py")
  run_script("denoise.py")  # Note: May refer to integrated denoising in Seperator.py
  run_script("speech_analyser.py")
  run_script("non-speech_analyser.py")

  # Extraction & Report
  speech_text = extract_speech()  # Clean regex
  nonspeech_text = extract_nonspeech() # Line parsing
  # Template fill and save
```

**Potential Limitations**: Brittle path assumptions; no parallel processing.

**Note on `denoise.py`** : Referenced in `main.py` but not provided—likely integrated into `Seperator.py` as `denoise_vocals()` .

# Performance and Deployment Notes

- **Dependencies**: Install via `pip install torch torchaudio demucs noisereduce openai-whisper tensorflow tensorflow-hub librosa pandas numpy` .

- **Runtime**: ~10-60s per minute of audio (CPU); scales with file length.

- **Accuracy Benchmarks** (estimated from model papers):

- Separation: SDR ~9dB (htdemucs).

- Transcription: WER <10% (Whisper large).

- Classification: mAP ~0.3 (YAMNet on AudioSet).

- **Extensions**: Add GPU support ( `device="cuda"` ), timestamps to transcription, or visualization (e.g., spectrograms).

# Conclusion

This pipeline delivers a robust, end-to-end solution for audio forensics, blending cutting-edge models like Demucs, Whisper, and YAMNet into a seamless workflow. It excels in scenarios like podcast analysis, surveillance audio review, or content moderation, providing actionable insights via a simple TXT report.

**Under the guidance of our team leader Harshita, we did this.**

Future enhancements could include real-time processing or multi-file batching. For implementation queries, refer to the provided scripts.