

PEOPLE'S DEMOCRATIC REPUBLIC OF ALGERIA  
MINISTRY OF HIGHER EDUCATION AND SCIENTIFIC RESEARCH  
SAAD DAHLEB BLIDA 01 UNIVERSITY  
DEPARTMENT OF COMPUTER SCIENCE



MASTER'S INTELLIGENT SYSTEMS ENGINEERING

---

**SPEECH PROCESSING**

REPORT

**A DELVE INTO VOICE ANALYSIS  
TECHNIQUES**

Made by :

ABDELATIF MEKRI

Academic year : 2023-2024

# Contents

I.	Abstract.....	4
II.	Generalities .....	4
1.	What is voice analysis.....	4
2.	Usage of voice analysis.....	4
3.	Differences between voice analytics, speech analytics, and text analytics .....	5
3.1.	Speech Analytics .....	6
3.2.	Voice Analytics .....	6
3.3.	Text Analytics .....	6
III.	Voice Analysis.....	6
1.	Voice generation and reception .....	6
2.	Mel Scale .....	8
3.	Methods of Audio analysis .....	9
3.1.	Fourier Transform .....	9
3.2.	Spectrogram.....	9
3.3.	Pitch Detection.....	10
3.4.	Mel-Frequency Cepstral Coefficients (MFCC) .....	11
3.5.	Waveform Analysis .....	11
3.6.	Audio Feature Extraction .....	12
3.7.	Machine Learning and Pattern Recognition .....	12
IV.	Implementation of Voice Analysis Methods.....	13
1.	Technical Implementation.....	13
1.1.	Data Handling and Audio Processing .....	13
1.2.	Machine Learning and Deep Learning: .....	13
1.3.	Speech Recognition and Text-to-Speech: .....	13
1.4.	Signal Processing and Feature Extraction: .....	13
1.5.	Visualization and Plotting: .....	14
V.	Conclusion.....	14
VI.	Refrences .....	15

## Table of Figures

Figure 1: Human speech communication (Holmes & Holmes, 2001).....	7
Figure 2: Source-channel model for a speech recognition system (Huang et al., 2001) .....	7
Figure 3: Mel Scale Representation .....	8
Figure 4: View of a signal in the time and frequency domain with Inverse Fourier Transform .....	9
Figure 5: A spectrogram plot example.....	10
Figure 6: Audio Signal to Pitch Estimation .....	10
Figure 7: Mel-Frequency Cepstral Coefficients (MFCC) .....	11
Figure 8: Waveform representation of a text speech.....	11
Figure 9: Typical pipeline used by audio deep learning models .....	12

## I. Abstract

Voice analysis techniques have emerged as a powerful tool across various domains, ranging from healthcare to security, psychology, and beyond. This report provides a comprehensive examination of the current landscape of voice analysis techniques, highlighting their efficacy and diverse applications. Beginning with an overview of the underlying principles and methodologies, the report delves into the advancements in voice analysis technology, including machine learning algorithms and signal processing methods.

Furthermore, the report explores the practical applications of voice analysis techniques, such as emotion recognition, lie detection, speech pathology diagnosis, and speaker identification. It discusses the challenges associated with voice data collection, processing, and interpretation, emphasizing the importance of accuracy, reliability, and ethical considerations.

Additionally, the report examines the potential impact of voice analysis techniques on society, including implications for privacy, security, and individual rights. Finally, it offers insights into future directions and areas of research, highlighting the opportunities for innovation and collaboration in this rapidly evolving field. Overall, this report provides valuable insights into the usage of voice analysis techniques and their significance in various domains.

## II. Generalities

### 1. What is voice analysis

Voice analysis refers to the use of specialized software to analyze spoken language, extracting insights and information from audio recordings. This process involves techniques such as speech recognition, natural language processing, sentiment analysis, and machine learning algorithms to interpret the content, tone, and context of conversations. Initially focused on transcribing speech into text, modern voice analytics tools now offer advanced capabilities such as conversation intelligence, sentiment analysis, speaker identification, insights generation, and predictive analytics. These tools are widely used across industries to improve customer experience, optimize processes, and drive business outcomes.<sup>1</sup>

### 2. Usage of voice analysis

Voice analysis offers several benefits across various domains:

- **Enhanced Customer Experience:** Voice analysis helps organizations understand customer sentiment, preferences, and pain points by analyzing interactions. This insight enables them to tailor services and solutions to meet customer needs effectively.
- **Improved Decision-Making:** By extracting actionable insights from voice data, organizations can make informed decisions regarding product development, marketing strategies, and customer service initiatives. Voice analysis provides valuable feedback that guides decision-making processes.

---

<sup>1</sup> The Team at CallMiner. (2024, February 28). What is voice analytics? Definition, tips and best practices. CallMiner.

- **Efficient Operations:** Voice analysis automates the process of analyzing spoken interactions, saving time and resources compared to manual review. This efficiency allows organizations to streamline operations and focus on strategic tasks rather than routine analysis.
- **Quality Assurance:** In customer service or sales environments, voice analysis helps monitor and evaluate interactions to ensure compliance with regulations, adherence to best practices, and consistency in messaging. It enables organizations to maintain high-quality standards across all customer touchpoints.
- **Personalized Services:** By understanding individual preferences and behaviors through voice analysis, organizations can personalize customer interactions and offerings. This personalization enhances customer satisfaction and fosters long-term loyalty.
- **Fraud Detection:** Voice analysis can identify suspicious patterns or anomalies in spoken interactions, helping organizations detect fraudulent activities such as identity theft or unauthorized access to sensitive information.
- **Healthcare Applications:** In healthcare, voice analysis techniques aid in diagnosing speech disorders, monitoring patient health remotely, and identifying early signs of cognitive decline or mental health issues.
- **Security Enhancement:** Voice analysis can be used for biometric authentication, verifying individuals based on unique vocal characteristics. This enhances security in applications such as access control, banking, and e-commerce.
- **Research and Insights:** Voice analysis provides researchers with valuable data for studying linguistic patterns, cultural trends, and social behaviors. It offers insights into communication dynamics and can contribute to academic studies and market research.
- **Continuous Improvement:** By analyzing voice data over time, organizations can identify trends, patterns, and areas for improvement. This continuous feedback loop enables iterative enhancements to products, services, and processes.

Overall, voice analysis offers significant benefits in terms of improving customer experience, operational efficiency, decision-making, security, and innovation across a wide range of industries and applications.

### 3. Differences between voice analytics, speech analytics, and text analytics

In the analysis of communication data, terms like speech analytics, voice analytics, and text analytics are often used interchangeably, blurring their distinct roles. However, upon closer examination, each term denotes a specific facet of communication analysis, offering unique insights into content, tone, and medium. This comparison seeks to elucidate the disparities between speech analytics, voice analytics, and text analytics, clarifying their individual focuses and functionalities. By delineating these differences, we aim to provide clarity on the array of tools available for analyzing communication data, empowering informed decision-making in their effective utilization.

### 3.1. Speech Analytics

**Definition:** Speech analytics involves identifying and examining the actual words spoken during a conversation.

**Medium:** It deals with spoken language, converting speech into text for analysis purposes.

**Purpose:** Its primary function is to analyze the content of spoken interactions, identifying patterns, trends, and frequently used phrases.

### 3.2. Voice Analytics

**Definition:** Voice analytics focuses on understanding how things are expressed in a conversation, such as tone, pitch, tempo, rhythm, and stress.

**Medium:** It analyzes the audio characteristics of speech without necessarily transcribing it into text.

**Purpose:** Its main goal is to use sentiment analysis to uncover the emotional aspects of communication, providing insights into the underlying intent and sentiment behind spoken words.

### 3.3. Text Analytics

**Definition:** Text analytics involves analyzing the meaning, intent, and emotional context of written language, encompassing text messages, emails, chats, and social media posts.

**Medium:** It deals exclusively with written or textual communication.

**Purpose:** Its objective is to extract insights from written content, including sentiment analysis, topic modeling, and entity recognition, to understand customer opinions, preferences, and behaviors.

In summary, while speech analytics focuses on the content of spoken conversations, voice analytics explores the emotional nuances of communication, and text analytics deals with written communication. Together, these technologies form the basis of conversation analytics, enabling organizations to gain a holistic understanding of customer interactions across various channels and mediums.

## III. Voice Analysis

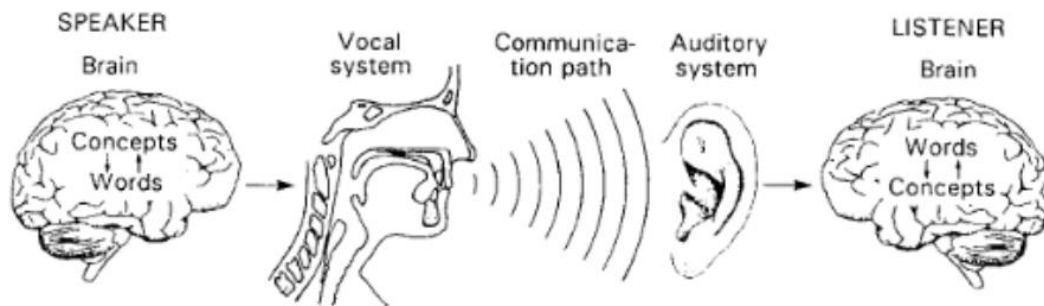
### 1. Voice generation and reception

Understanding how speech is produced and perceived is essential for delving into speech recognition. This process serves as the bridge between speakers and listeners, covering the journey from creating speech sounds to interpreting them.

Speech production starts with airflow from the lungs passing through the vocal cords, which vibrate to produce sound. These sounds are then shaped by the tongue, lips, and palate into recognizable speech sounds. Finally, they're emitted as acoustic signals.

On the listener's side, speech perception involves capturing and processing these acoustic signals through the auditory system. The ear captures the signals and converts them into electrical impulses sent to the brain. The brain decodes and interprets these signals, recognizing sound patterns, phonemes, and words to understand the spoken language.

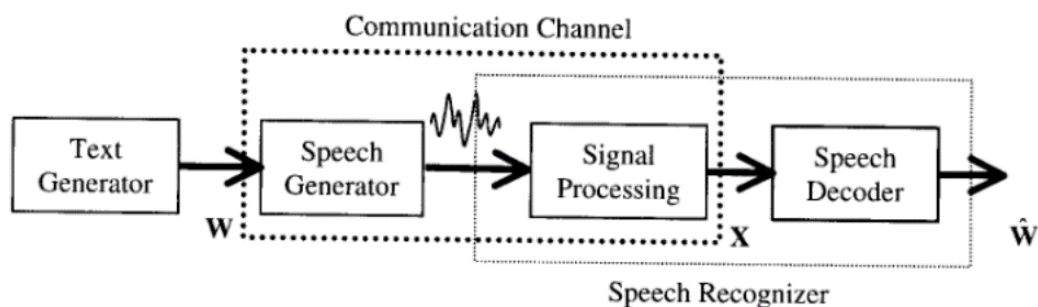
This interaction between speech production and perception forms the foundation of speech recognition technology. By understanding these processes, researchers can develop algorithms and systems to accurately recognize and understand human speech.



2

Figure 1: Human speech communication (Holmes & Holmes, 2001)

Speech recognition systems aim to mimic the human speech communication system. Figure 2 illustrates a source-channel model commonly used in speech recognition systems.



3

Figure 2: Source-channel model for a speech recognition system (Huang et al., 2001)

Human speech communication serves the purpose of transferring ideas. Initially, these ideas originate within the speaker's brain and are then formulated into a sequence of words, denoted as the source word sequence  $\hat{W}$ . This sequence is transformed by the speaker's vocal system, modeled as the speech generator component, into a speech signal waveform. This waveform travels through the air, acting as a noisy communication channel, and may be subject to interference from external noise sources.

Upon reaching the listener, the acoustical signal is perceived by the human auditory system. The listener's brain processes this waveform to comprehend its content, completing the communication process. This perception process is akin to the signal processing and speech decoding components of

<sup>2</sup> Alcaraz Meseguer, N. (2009). Speech Analysis for Automatic Speech Recognition (Master's thesis). Norwegian University of Science and Technology.

a speech recognizer. Their objective is to process and decode the acoustic signal  $X$  into a word sequence  $\hat{W}$ , ideally resembling the original word sequence  $\hat{W}$ .

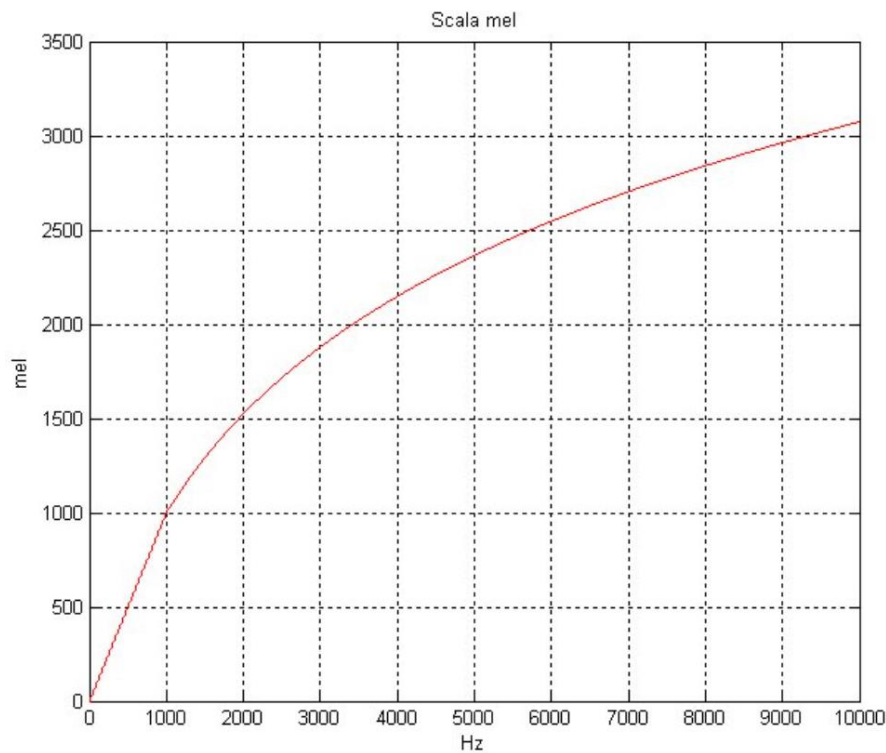
## 2. Mel Scale

The Mel scale is a perceptual scale of pitches, which approximates the human ear's response to different frequencies of sound. It is named after the researcher, Stevens, who conducted early experiments in auditory perception. The Mel scale is particularly relevant in audio signal processing, speech recognition, and music technology.

Unlike the linear scale of frequency in Hertz (Hz), where each step corresponds to an equal increase in pitch, the Mel scale is nonlinear and follows the perceived pitch changes in human hearing. It is based on the idea that the ear's response to changes in frequency is not linear but rather logarithmic.

The Mel scale is commonly used for tasks like feature extraction in speech recognition, where it helps capture the characteristics of human speech more accurately. This scale is divided into critical bands, each representing a range of frequencies that the human ear perceives as one unit.

Mel scale provides a more accurate representation of how humans perceive pitch, making it valuable in various audio-related applications.



4

Figure 3: Mel Scale Representation

---

<sup>4</sup> Karam, S., Dutta, I., Ijaz, M. F., & Singh, P. K. (2021). TLEFuzzyNet: Fuzzy Rank-Based Ensemble of Transfer Learning Models for Emotion Recognition From Human Speeches.

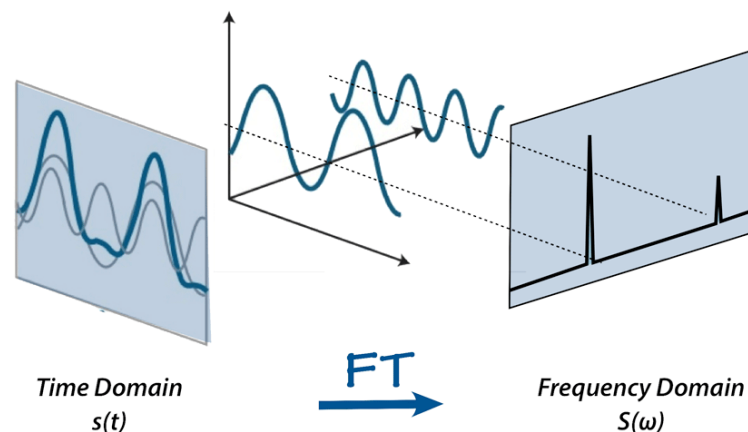


### 3. Methods of Audio analysis

Audio analysis encompasses various methods and techniques for examining and extracting meaningful information from audio signals. It plays a crucial role in tasks such as signal processing, feature extraction, classification, and recognition across a wide range of applications.

#### 3.1. Fourier Transform

This mathematical operation decomposes a time-domain signal into its constituent frequencies, providing insights into the frequency content of an audio signal.



5

Figure 4: View of a signal in the time and frequency domain with Inverse Fourier Transform

#### 3.2. Spectrogram

A visual representation of the frequency content of a signal over time, useful for analyzing how frequency changes over time and detecting transient events.

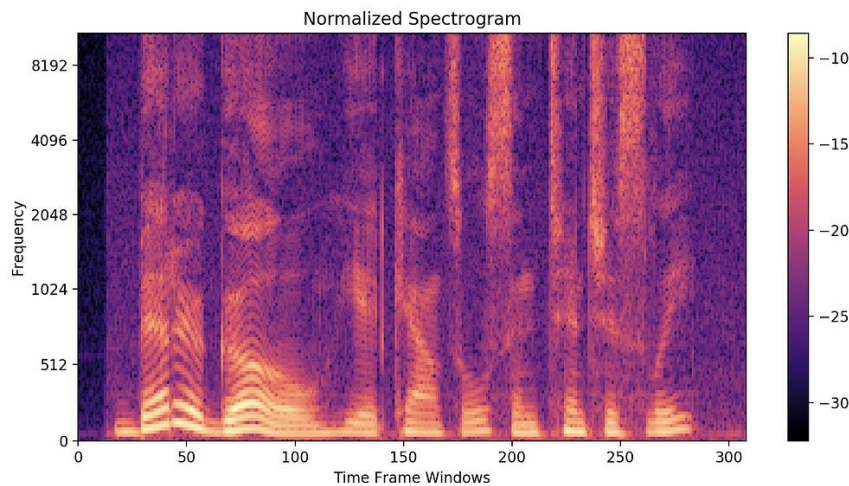
In a speech recognition task, preserving both frequency and time information is crucial for accurate prediction of spoken words. While the Fourier Transform (FFT) provides frequency values, it lacks temporal context. Spectrograms offer a solution by visually representing frequencies over time.

A spectrogram displays frequency content over time, with time on one axis, frequency on the other, and magnitude (amplitude) represented by colors. Bright colors indicate strong frequencies, with smaller frequencies typically appearing brighter. This representation allows us to retain the time information lost in the FFT plot while capturing the frequency content of the audio signal.

By analyzing spectrograms, our recognition system can extract features that incorporate both frequency and time information. This enables the system to recognize spoken words based on their spectral characteristics over time, improving accuracy and performance in speech recognition tasks.

---

<sup>5</sup> Chaudhary, K. (2020, January 19). Understanding audio data, Fourier transform, FFT and spectrogram features for a speech recognition system.

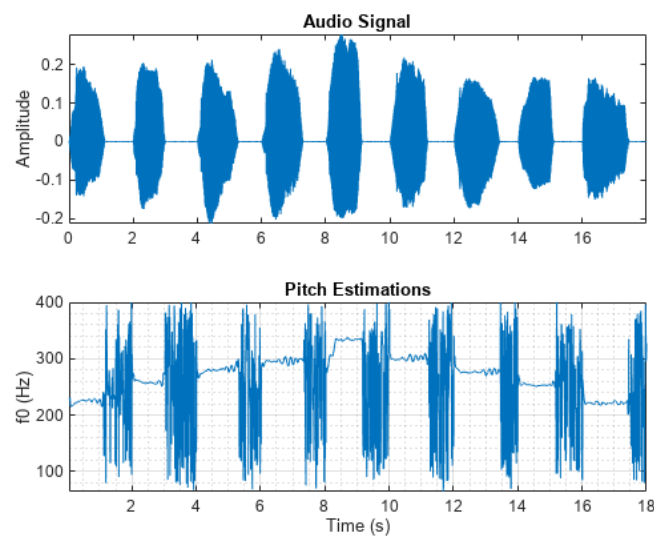


6

Figure 5: A spectrogram plot example

### 3.3. Pitch Detection

Also known as fundamental frequency estimation, is a vital component in various audio applications, particularly in speech and music analysis. It involves determining the fundamental frequency (or pitch) of a periodic signal, which corresponds to the perceived pitch of the sound. Here's an overview of pitch detection:



7

Figure 6: Audio Signal to Pitch Estimation

The fundamental frequency represents the rate at which the waveform repeats, determining the perceived pitch of the sound. Pitch detection is crucial for tasks such as music transcription, speech analysis, and voice recognition or where accurate estimation of pitch is required, such as emotion detection, speaker identification, and intonation analysis.

<sup>6</sup> Chaudhary, K. (2020, January 19). Understanding audio data, Fourier transform, FFT and spectrogram features for a speech recognition system.

<sup>7</sup> MathWorks. (n.d.). pitch. Retrieved May 11, 2024, from <https://www.mathworks.com/help/audio/ref/pitch.html>

### 3.4. Mel-Frequency Cepstral Coefficients (MFCC)

A feature extraction technique widely used in speech and audio analysis, designed to match human auditory perception, by capturing key spectral features of the signal, such as the shape of the vocal tract and the distribution of energy across different frequency bands.

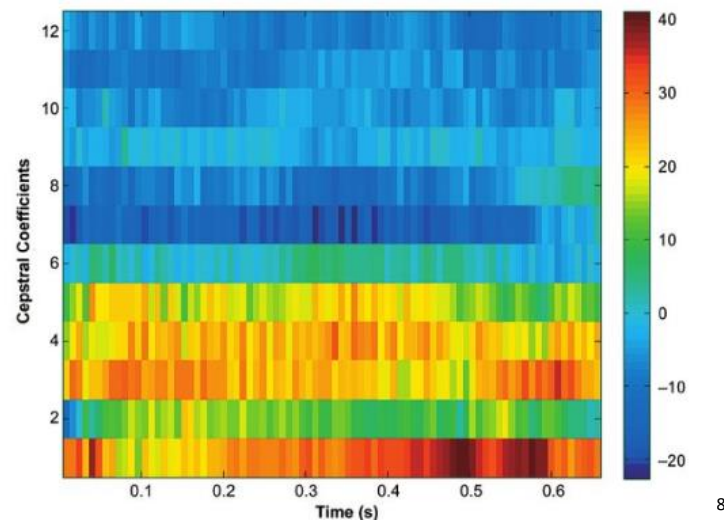


Figure 7: Mel-Frequency Cepstral Coefficients (MFCC)

MFCCs find applications in audio content analysis tasks such as emotion detection, language identification, and sound classification.

### 3.5. Waveform Analysis

It provides insights into how the signal varies over time and can reveal important features such as the presence of specific sounds, changes in intensity, and the overall structure of the audio.

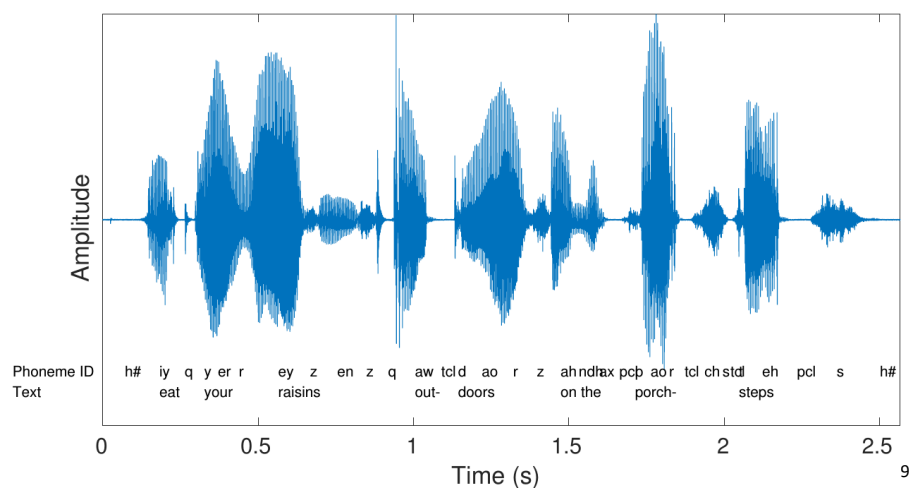


Figure 8: Waveform representation of a text speech

<sup>8</sup> Engel, Z. W., Kłaczyński, M., & Wszolek, W. (2007). A Vibroacoustic Model of Selected Human Larynx Diseases.

<sup>9</sup> Bäckström, T., Räsänen, O., Zewoudie, A., Pérez Zarazaga, P., Koivusalo, L., Das, S., Gómez Mellado, E., Bouafif Mansali, M., & Ramos, D. (n.d.). Waveform representations.

Waveform analysis is useful for tasks where temporal characteristics of the audio signal are important, such as event detection, onset detection, and segmentation or to manipulate audio waveforms to perform tasks like cutting, trimming, and aligning audio clips.

### 3.6. Audio Feature Extraction

involves capturing specific characteristics or attributes of an audio signal that are relevant to the task at hand. These features serve as inputs to various audio analysis algorithms and machine learning models, enabling tasks such as classification, clustering, and regression.

### 3.7. Machine Learning and Pattern Recognition

These methods enable the development of models that can automatically learn patterns and relationships from audio data, leading to advancements in tasks such as audio classification, speech recognition, and sound event detection. Here's an overview of machine learning and pattern recognition in the context of audio analysis

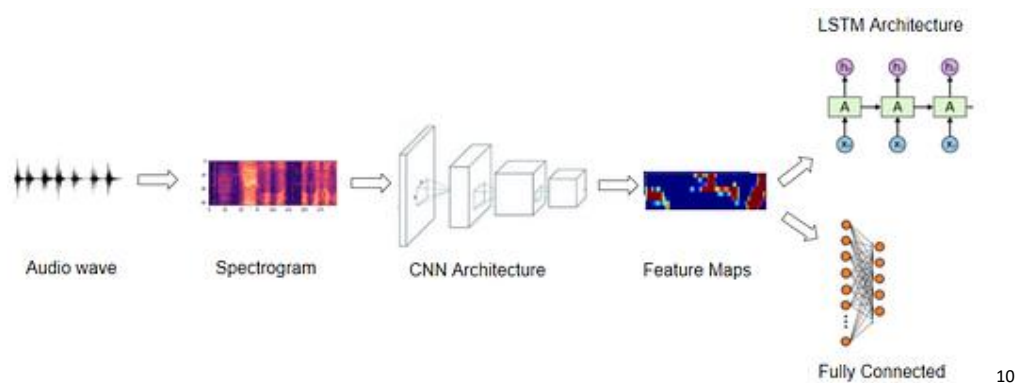


Figure 9: Typical pipeline used by audio deep learning models

Machine learning and pattern recognition techniques are essential in tasks where manual analysis or rule-based approaches are not feasible due to the complexity or variability of the audio data.

Especially where the goal is to categorize audio signals into predefined classes or to build models to transcribe spoken language into text, enabling applications such as virtual assistants, voice-controlled devices, and dictation software.

---

<sup>10</sup> Doshi, K. (2021, February 11). Audio deep learning made simple (Part 1): State-of-the-art techniques. Towards Data Science.

## IV. Implementation of Voice Analysis Methods

### 1. Technical Implementation

Implementing models using voice analysis methods involves several steps, including data preprocessing, feature extraction, model training, and evaluation . All of these steps require some powerful libraries to handle the specific nature of each task .

#### 1.1. Data Handling and Audio Processing

- Librosa: A Python package for music and audio analysis. It provides tools for loading audio files, extracting features (e.g., MFCCs, spectrograms), and performing various audio processing tasks.
- PyAudio: A Python library for audio I/O. It enables recording and playback of audio data from microphones and speakers, which is useful for collecting audio samples for training datasets.
- FFmpeg: A multimedia framework that can be used for reading, writing, and processing audio and video files. It provides powerful features for audio manipulation and conversion.

#### 1.2. Machine Learning and Deep Learning:

- Scikit-learn: A comprehensive machine learning library for Python. It includes a wide range of algorithms for classification, regression, clustering, and dimensionality reduction, as well as tools for model evaluation and validation.
- TensorFlow: An open-source deep learning framework developed by Google. It provides tools for building and training neural networks, including high-level APIs like Keras for easier model construction and training.
- PyTorch: A deep learning framework developed by Facebook. It offers dynamic computation graphs and a flexible, imperative programming interface, making it popular for research and development of neural network models.

#### 1.3. Speech Recognition and Text-to-Speech:

- SpeechRecognition: A Python library that provides easy access to several speech recognition APIs (such as Google Speech Recognition, IBM Watson Speech to Text, and CMU Sphinx). It enables the integration of speech recognition capabilities into Python applications.
- pyttsx3: A Python library for text-to-speech conversion. It allows developers to synthesize natural-sounding speech from text using different speech engines.

#### 1.4. Signal Processing and Feature Extraction:

- NumPy: A fundamental package for scientific computing with Python. It provides support for multi-dimensional arrays and matrices, essential for numerical operations and data manipulation.
- SciPy: A library for scientific computing in Python. It builds on NumPy and provides additional functionality for signal processing, interpolation, optimization, and more.

- Pandas: A powerful data analysis library for Python. It offers data structures and functions for handling structured data and time-series data, which can be useful for organizing and analyzing audio features extracted from datasets.

#### 1.5. Visualization and Plotting:

- Matplotlib: A plotting library for Python. It provides a MATLAB-like interface for creating static, interactive, and animated visualizations of data, including plots of audio waveforms, spectrograms, and feature distributions.
- Seaborn: A statistical data visualization library based on Matplotlib. It provides a high-level interface for drawing informative and attractive statistical graphics, which can be useful for visualizing relationships between audio features and target labels.

## V. Conclusion

Voice analysis has emerged as a critical field within the broader domain of audio signal processing, offering profound insights into human communication, interaction, and behavior. Through the application of sophisticated techniques such as speech recognition, emotion detection, and speaker identification, voice analysis has revolutionized numerous industries and applications, spanning from telecommunications and customer service to healthcare and entertainment.

By harnessing the power of machine learning, deep learning, and signal processing algorithms, voice analysis enables the extraction of valuable information from audio data, ranging from identifying spoken words and phrases to discerning underlying emotions and sentiments. This wealth of information has paved the way for the development of innovative solutions, including virtual assistants, speech-to-text transcription services, personalized recommendations, and biometric authentication systems.

Moreover, voice analysis holds immense promise for addressing societal challenges and enhancing human-computer interaction. For instance, in healthcare, voice-based diagnostics and monitoring systems offer non-invasive means of assessing patients' health conditions and emotional states. In education, speech recognition technologies facilitate language learning and literacy improvement by providing real-time feedback and personalized tutoring.

Despite its transformative potential, voice analysis also poses significant technical and ethical considerations. Challenges such as noise robustness, privacy protection, and bias mitigation require ongoing research and development efforts to ensure the responsible and equitable deployment of voice analysis technologies.

In conclusion, voice analysis stands at the forefront of technological innovation, driving advancements across diverse domains and revolutionizing the way we interact with machines and each other. With continued advancements in algorithm development, data collection, and interdisciplinary collaboration, voice analysis holds the promise of unlocking new frontiers in communication, understanding, and human-machine symbiosis.

## VI. References

### PAPERS & ARTICLES

- [1] Alcaraz Meseguer, N. (2009). Speech Analysis for Automatic Speech Recognition (Master's thesis). Norwegian University of Science and Technology, Department of Electronics and Telecommunications.
- [2] Karam, S., Dutta, I., Ijaz, M. F., & Singh, P. K. (2021). TLEFuzzyNet: Fuzzy Rank-Based Ensemble of Transfer Learning Models for Emotion Recognition From Human Speeches. IEEE Access, 99, 1-1. <https://doi.org/10.1109/ACCESS.2021.3135658>

### WEBSITES

- [3] Chaudhary, K. (2020, January 19). Understanding audio data, Fourier transform, FFT and spectrogram features for a speech recognition system: An introduction to audio data analysis (sound analysis) using Python. Towards Data Science. <https://towardsdatascience.com/understanding-audio-data-fourier-transform-fft-spectrogram-and-speech-recognition-a4072d228520>
- [4] The Team at CallMiner. (2024, February 28). What is voice analytics? Definition, tips and best practices. CallMiner. <https://callminer.com/blog/what-is-voice-analytics-definition-tips-best-practices-and-challenges-of-voice-analytics>
- [5] [https://www.researchgate.net/publication/357036500\\_TLEFuzzyNet\\_Fuzzy\\_Rank-Based\\_Ensemble\\_of\\_Transfer\\_Learning\\_Models\\_for\\_Emotion\\_Recognition\\_From\\_Human\\_Speeches/figures?lo=1](https://www.researchgate.net/publication/357036500_TLEFuzzyNet_Fuzzy_Rank-Based_Ensemble_of_Transfer_Learning_Models_for_Emotion_Recognition_From_Human_Speeches/figures?lo=1)
- [6] MathWorks. (n.d.). pitch. Retrieved May 11, 2024, from <https://www.mathworks.com/help/audio/ref/pitch.html>
- [7] Engel, Z. W., Kłaczyński, M., & Wszolek, W. (2007). A Vibroacoustic Model of Selected Human Larynx Diseases. International Journal of Occupational Safety and Ergonomics (JOSE), 13(4), 367-379. <https://doi.org/10.1080/10803548.2007.11105094>
- [8] Bäckström, T., Räsänen, O., Zewoudie, A., Pérez Zarazaga, P., Koivusalo, L., Das, S., Gómez Mellado, E., Bouafif Mansali, M., & Ramos, D. (n.d.). Waveform representations. SpeechProcessingBook. Retrieved from <https://speechprocessingbook.aalto.fi/Representations/Waveform.html>
- [9] Doshi, K. (2021, February 11). Audio deep learning made simple (Part 1): State-of-the-art techniques. Towards Data Science. <https://towardsdatascience.com/audio-deep-learning-made-simple-part-1-state-of-the-art-techniques-da1d3dff2504>