

Speaker Verification in Real-World Applications: Advances and Emerging Trends

N. Lohith Sai
Department of Computer science and
Engineering,
Koneru Lakshmaiah
Education Foundation,
Vaddeswaram, Andhra Pradesh, India

K. Kishore
Department of Computer science and
Engineering,
Koneru Lakshmaiah
Education Foundation,
Vaddeswaram, Andhra Pradesh, India

Y. Rupesh Kumar Reddy
Department of Computer science and
Engineering,
Koneru Lakshmaiah
Education Foundation,
Vaddeswaram, Andhra Pradesh, India

K. Gowtham Sai
Department of Computer science and
Engineering,
Koneru Lakshmaiah
Education Foundation,
Vaddeswaram, Andhra Pradesh, India

ABSTRACT

We outline the main components of the MIT Lincoln Laboratory's Gaussian mixture model (GMM)-based speaker verification system in this work. This system has been used to successfully verify speakers in many NIST Speaker Recognition Evaluations (SREs). The approach is based on the likelihood ratio test for verification and employs universal background models (UBMs) for alternative speaker representation, basic but to construct speaker models from the UBM, functional GMMs for chances distributions and a sort of Bayesian adaptation are used. Furthermore, the development and use of a handset detector, as well as score normalisation, which significantly enhance verification performance, are detailed and addressed. The presentation concludes with realistic performance benchmarks and experiments on system behaviour using NIST SRE corpora. The research directions being jointly pursued at the Speech and Vision Laboratory of the Indian Institute of Technology Madras and the Anthropoc Signal Processing Group of the Oregon Graduate Institute are discussed in this paper. The current approaches to speaker verification rely on Gaussian mixture models (GMM) to characterise the speaker's characteristics. If target speakers use a different phone handset from that used during training, the performance of these systems suffers noticeably. Both utterance-based mean subtraction (MS) and relative spectrum (RASTA) filtering are common techniques for channel normalisation. In this paper, we describe a revolutionary method to filter design that can normalise the variability introduced by various phone handsets. The estimated second-order statistics of the handset are used to design the filter. The process of authenticating a speaker's claim as truthful or untrue entails evaluating the voice signal. Deep neural networks are one of the most successful. Implementation of complicated nonlinear models to learn unique and invariant data characteristics. They have been utilised in voice recognition tasks and have demonstrated the potential to be used for speaker recognition as well. In this paper, we analyse and discuss Deep Neural Network (DNN)

approaches utilised in speech verification Systems

I. INTRODUCTION

Speaker verification, an essential biometric technique, aims to verify or authenticate individuals based on their unique voice characteristics. It involves comparing the voice of a claimed speaker with a stored voiceprint to establish their identity. Gaussian Mixture Models (GMMs) have gained popularity as a prominent approach for speaker verification. GMMs are powerful statistical models capable of accurately representing complex data distributions, making them suitable for modeling voice features. This paper delves into the process of speaker verification using GMM models and emphasizes their significance in this domain.

1.1 Motivation:

Speaker verification has gained significant attention in the field of biometric authentication due to its potential for secure and convenient user identification. Traditional methods of authentication, such as passwords or PINs, are vulnerable to various security threats, including unauthorized access and identity theft. Speaker verification offers an alternative approach by leveraging the unique vocal characteristics of individuals. Motivated by the need for reliable and user-friendly authentication systems, this project aims to explore and develop an effective speaker verification system.

1.2 Background:

Speaker verification is a biometric technology that focuses on verifying the claimed identity of an individual based on their voice. Each person has unique vocal characteristics, including pitch, tone, speech patterns, and pronunciation, which can be used to establish their identity. The field of speaker verification has seen significant advancements in recent years, thanks to the progress in speech signal processing, machine learning, and deep learning techniques. These advancements have paved the way for more accurate and robust speaker verification systems.

1.3 Objectives:

The primary objective of this speaker verification project is to develop a reliable and accurate system that can accurately verify the identity of individuals based on their speech. The specific objectives include:

Investigating various feature extraction techniques, such as MFCCs, LPC, PLP, and deep learning-based features, to identify the most effective approach for capturing discriminative speaker information.

Exploring different modeling approaches, including Gaussian Mixture Models (GMMs), Support Vector Machines (SVMs), Deep Neural Networks (DNNs), Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and speaker embeddings, to build robust and accurate speaker models.

Evaluating and comparing different decision-making techniques, such as score normalization, threshold selection, Bayesian approaches, and fusion of multiple models, to determine the optimal strategy for making reliable verification decisions.

Addressing the challenges faced in speaker verification, such as data availability and collection, variability in speech signals, channel and environmental conditions, impostor detection, and computational complexity, to enhance the system's performance and robustness.

Conducting rigorous evaluations of effectiveness on the built speaker verification technology using standard evaluation metrics and benchmark datasets to measure its efficacy and compare it with current cutting-edge methods.

By achieving these objectives, this project aims to contribute to the advancement of speaker verification technology and they offer useful insights into the design and implementation of trustworthy and effective verification of speakers systems. II. Speaker Verification using GMM Models

2.1 Definition and Purpose:

Speaker verification, also known as speaker authentication or voice verification, is a biometric technology that aims to authenticate the claimed identity of an individual based on their unique voice characteristics. The goal of speaker verification is to assess if the speech sample supplied by anyone fits the voice profile linked with their stated identity.

Speaker verification systems are designed to address the need for secure and reliable user authentication in various applications. They find applications in voice-controlled systems, access control, phone banking, forensic investigations, and other scenarios where voice is used as a means of identity verification.

2.2 Key Components:

A typical speaker verification system consists of the following key components:

a. Enrollment: Throughout the enrolling step, the system collects the user's speech sample and extracts important data features to create a speaker model or template. This template represents the unique characteristics of the individual's voice and serves as a reference for future verification.

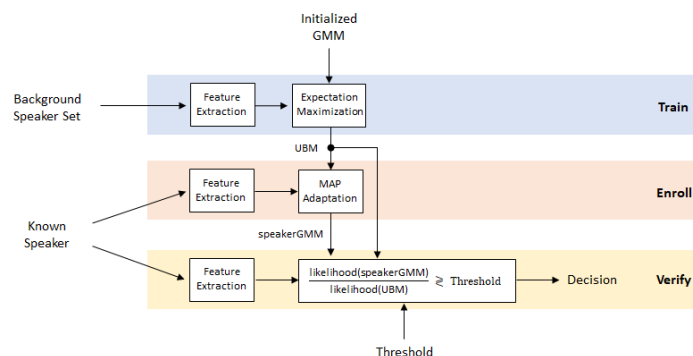
b. Feature Extraction: In the feature extraction stage, acoustic features are extracted from the voice signal. These features capture the distinguishing characteristics of the individual's voice, such as Mel-Frequency Cepstral Coefficients (MFCCs), Linear Predictive Coding (LPC) coefficients, or deep learning-based representations.

c. Speaker Modeling: Speaker modeling involves creating a

mathematical representation of the speaker based on the extracted features. This can be done using techniques such as Gaussian Mixture Models (GMMs), Support Vector Machines (SVMs), Deep Neural Networks (DNNs), or speaker embeddings.

d. Verification: In the verification phase, the system compares the voice sample provided during the verification attempt with the enrolled speaker model. The system calculates a similarity score or distance metric between the two voice samples to determine whether they belong to the same speaker or not.

e. Decision-Making: Based on the similarity score or distance metric, a decision is made to accept or reject the claimed identity. This choice is often made by using score normalisation procedures and threshold selection, or using Bayesian approaches to ensure a reliable verification decision.



2.3 Performance Evaluation Metrics:

Several tests are used to evaluate the performance of a speaker verification system. evaluation metrics are commonly used:

a. Equal Error Rate (EER): EER denotes the point at which the false acceptance rate (FAR) and false rejection rate (FRR) are equal. It offers a measure of the system's overall performance in accurately accepting legitimate speakers, while rejecting impostors.

b. False Acceptance Rate (FAR): FAR indicates the proportion of imposter efforts that are wrongly identified as actual speakers by the system. A lower FAR indicates better system security against impostor attacks.

c. False Rejection Rate (FRR): FRR reflects the proportion of valid speaker efforts that are mistakenly rejected by the system. A lower FRR indicates better system usability and user convenience.

d. Receiver Operating Characteristic (ROC) Curve: The ROC curve plots the system's performance by varying the threshold for decision-making. It shows the trade-off between the FAR and FRR and provides insights into the system's performance at different operating points.

e. Detection Cost Function (DCF): DCF combines the FAR and FRR with associated costs to provide an overall measure of system performance that incorporates both security and usability considerations.

f. Accuracy: Accuracy indicates the overall accuracy of the system's verification judgements. considering both true positives and true negatives.

III Feature Extraction Techniques

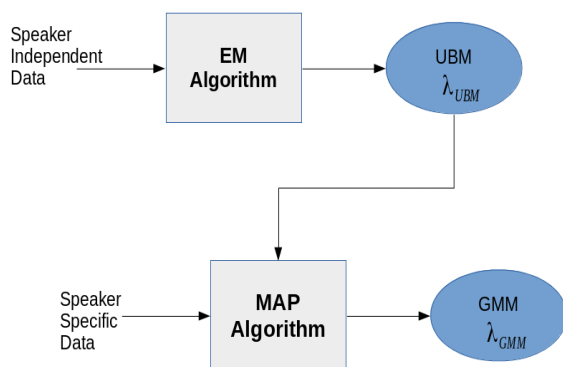
3.1 Mel-Frequency Cepstral Coefficients (MFCCs):

In speaker verification systems, Mel-Frequency Cepstral Coefficients (MFCCs) are a popular feature extraction approach. Several stages are involved in the extraction of MFCC. First, the voice signal is separated into small frames of 20-40 milliseconds. The range of colours is then balanced by pre-emphasizing each frame. To limit spectral leaks, every image is then subjected to a window function (e.g., Hamming window). The value of the spectra is obtained by

computing the Fourier transform. The mel-filterbank is then applied to the magnitude spectrum, capturing its spectral energy distribution using a collection of triangle filters placed on a mel scale. The filterbank outputs exponential is used to approximate the human auditory sense of loudness.

3.2 Linear Predictive Coding (LPC):

Linear Predictive Coding (LPC) is another feature extraction technique commonly employed in speaker verification. LPC represents the speech signal as a linear blend of previous samples, where the coefficients are estimated using methods like autocorrelation or the Yule-Walker equation. The LPC analysis aims to estimate the vocal tract resonances and formant frequencies, which are essential for speaker identification. The LPC coefficients capture the spectral envelope of the speech signal and provide a compact representation of the vocal tract characteristics. LPC features have been widely used in speaker verification systems, especially in applications where capturing the vocal tract information is crucial.



3.3 Perceptual Linear Prediction (PLP):

Perceptual Linear Prediction (PLP) result is an extension of LPC that incorporates psychoacoustic principles to improve the perceptual relevance of the extracted features. PLP takes into account the properties of human hearing, such as frequency resolution and masking effects, to obtain features that align better with human perception. PLP is concerned with the non-linear frequency spacing of crucial bands in the human hearing system and applies auditory masking models to emphasize perceptually relevant features. By considering these psychoacoustic principles, PLP coefficients capture important perceptual information in speech signals. PLP has been found to improve the effectiveness of speaker verification systems by incorporating human auditory perception.

3.4 Deep Learning-based Features:

Deep learning has shown great promise in speaker verification tasks. Deep learning-based feature extraction methods leverage neural networks, such as convolutional neural networks (CNNs) or recurrent neural networks (RNNs), to learn discriminative representations directly from raw speech signals. These models can automatically capture hierarchical and contextual information, allowing for more robust and context-aware speaker representations. Deep learning-based features, also known as speaker embeddings, have achieved state-of-the-art performance in speaker verification and have been widely adopted in many commercial systems. Speaker embeddings are low-dimensional representations learned by deep neural networks, which capture speaker-specific information. They provide a compact and discriminative feature representation for speaker verification.

In a speaker verification project, these feature extraction techniques can be explored and evaluated to determine which approach yields the best performance for the specific dataset and application. It is common to experiment with different combinations of feature extraction techniques and evaluate their impact on system performance. Additionally, hybrid approaches that combine multiple

feature extraction techniques or incorporate domain-specific knowledge can also be explored to increase the speaker verification system's accuracy and resilience.

IV Speaker Modelling Approaches:

4.1 Gaussian Mixture Models (GMMs):

Gaussian Mixture Models (GMMs) have been widely used in speaker verification systems. GMMs represent the distribution of speech features in a high-dimensional space. In speaker modeling, GMMs are used to model the feature space of individual speakers by estimating the parameters of multiple Gaussian components. The chance distribution of the speaker's characteristics is represented by each Gaussian component. During verification, the GMMs are used to calculate the likelihood of the test utterance given the speaker model. GMM-based speaker modeling has been successful in many speaker verification applications, especially with the use of Maximum Likelihood Linear Transform (MLLT) or Joint Factor Analysis (JFA) to further improve modeling accuracy.

4.2 Support Vector Machines (SVMs):

Support Vector Machines (SVMs) are another popular approach for speaker modeling in speaker verification systems. SVMs are binary classifiers that learn a decision boundary in a high-dimensional feature space. In speaker verification, SVMs can be trained to classify whether a given voice sample belongs to the target speaker or not. The SVMs learn a hyperplane that maximally separates the positive and negative examples in the feature space. SVMs have been effective in speaker verification, particularly when combined with appropriate kernel functions to handle non-linear separability. Commonly used kernels in speaker verification include linear, polynomial, and radial basis function (RBF) kernels.

4.3 Deep Neural Networks (DNNs):

Deep Neural Networks (DNNs) have revolutionized various domains, including speaker verification. DNNs are composed of multiple layers of interconnected artificial neurons that can learn complex representations from input data. In speaker modeling, DNNs can be used to learn discriminative speaker embeddings directly from the raw speech signals. DNNs, particularly architectures like Feedforward Neural Networks (FNNs) or Multi-Layer Perceptrons (MLPs), have been applied successfully to learn speaker-specific features that capture the unique characteristics of individual speakers. DNN-based approaches have shown superior performance in speaker verification tasks, especially when trained on large-scale datasets.

4.4 Convolutional Neural Networks (CNNs):

Convolutional Neural Networks (CNNs) have primarily been associated with image processing tasks but have also shown promise in speaker verification. CNNs can capture local patterns and spatial dependencies in the speech spectrogram, which is particularly useful for capturing spectral characteristics. By applying convolutional layers with filters of different sizes and pooling operations, CNNs can extract hierarchical and discriminative features from speech data. CNNs have been used for speaker verification by either processing the speech spectrogram directly or extracting features from intermediate layers and feeding them into subsequent modeling techniques.

4.5 Recurrent Neural Networks (RNNs):

Recurrent Neural Networks (RNNs) are well-suited for capturing temporal dependencies in sequential data, making them a valuable tool for speaker verification. RNNs, such as Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU), can model sequential dependencies in the time domain. Consequently, they are useful in modelling voice signals. RNNs can analyse variable-length voice sequences and collect long-term contextual information. They have been used for speaker verification by taking sequences of speech features as input and learning speaker-specific representations that encapsulate temporal dynamics. RNNs can be combined with other modeling approaches, such as GMMs or DNNs, for improved

performance.

V Decision Making Techniques:

5.1 Score Normalization:

Score normalization is a crucial step in speaker verification systems to ensure robust and reliable decision-making. In speaker verification, the similarity scores or distances obtained from the speaker modeling phase may vary across different speakers, utterances, or recording conditions. Score normalization approaches seek to normalise scores in order to make them more similar and consistent. Z-score normalisation is a frequent score normalisation approach. It transforms the similarity scores into standard Z-scores by subtracting the mean and dividing by the standard deviation of the scores obtained from a reference population or a development set. Z-score normalization helps in aligning the scores and making them comparable across different trials. Other score normalization techniques include T-norm, where the scores are linearly transformed to a target range, and histogram equalization, where the score distribution is matched to a predefined target distribution. Normalizing the scores reduces the influence of inter-speaker and inter-session variabilities, making the system more robust and reliable.

5.2 Threshold Selection:

Threshold selection is an important decision-making process in speaker verification systems. It involves determining a decision threshold that discriminates between genuine speakers and impostors based on the similarity scores or distances calculated during verification. Threshold selection can be performed using various methods. One common approach is to set a fixed threshold based on a predefined operating point that balances the false acceptance rate (FAR) and false rejection rate (FRR). The equal error rate (EER) threshold, which reflects the point at which FAR and FRR are equal, is frequently used as a guideline for thresholds selection. Another method is to create adaptive or dynamic thresholds based on the application's particular needs or the intended trade-off between safety and convenience. This can be accomplished by strategies such as cost-based thresholding, which takes into account the costs of mistaken acceptance and false rejection., or using machine learning techniques to learn the optimal threshold from training data.

5.3 Bayesian Approaches:

Bayesian approaches have been widely used in speaker verification to incorporate prior knowledge and statistical modeling into the decision-making process. Bayesian techniques provide a principled framework for combining evidence from multiple sources and making decisions based on probabilistic reasoning. One common Bayesian approach is the Bayesian likelihood ratio test, which calculates the likelihood ratio between the hypotheses of the target speaker and the impostor. The likelihood ratio is compared against a threshold to make the final decision.

Another Bayesian technique is the Bayesian belief network, which models the dependencies among different sources of evidence and combines them to compute the posterior probability of the target speaker's identity. Bayesian networks allow for the integration of multiple features, models, and contextual information in a unified framework.

5.4 Fusion of Multiple Models:

Fusion of multiple models is a powerful technique to improve the performance of speaker verification systems. It involves combining the decisions or similarity scores from multiple individual models to make a final decision. Model fusion can exploit complementary information from different models and enhance the overall system's accuracy and robustness. noise, channel effects, and emotional states.

VI Challenges in Speaker Verification:

6.1 Data Availability and Collection:

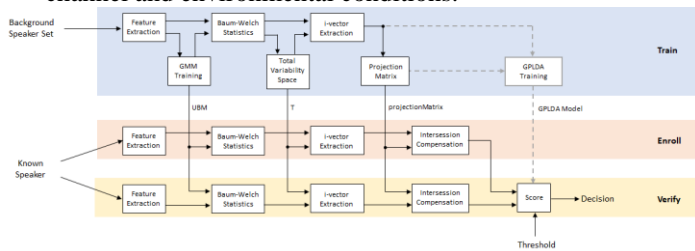
One of the key challenges in speaker verification is the availability and collection of sufficient and diverse training data. Building accurate and robust speaker models requires a large amount of labeled data from different speakers, speaking styles, languages, and demographics. However, collecting such data can be time-consuming, expensive, and challenging, especially for specific target populations or rare languages. Limited data can lead to overfitting or generalization issues, impacting the performance of the speaker verification system. Addressing data availability challenges may involve data augmentation techniques, data synthesis, or leveraging transfer learning from related tasks.

6.2 Variability in Speech Signals:

Speech signals are highly variable due to factors such as speaker-dependent characteristics, emotional state, accent, speaking rate, and coarticulation. This variability poses a significant challenge in speaker verification systems, as models need to generalize well across different variations while being sensitive to speaker-specific information. Addressing variability requires robust feature extraction techniques, robust modeling approaches, and adaptation mechanisms to capture and model speaker-specific information while mitigating the effects of non-speaker-related variations.

6.3 Channel and Environmental Conditions:

Speaker verification systems are often deployed in real-world scenarios where speech signals are acquired through various channels and under different environmental conditions. Variations in microphone quality, noise levels, reverberation, and transmission artifacts can degrade the performance of the system. Robustness to channel and environmental conditions is crucial to ensure reliable performance across different deployment scenarios. Techniques such as channel compensation, noise reduction, and robust feature extraction can help address the challenges associated with varying channel and environmental conditions.



6.4 Impostor Detection and Spoofing Attacks:

One of the critical challenges in speaker verification is the detection of impostors and spoofing attacks. Impostors may attempt to mimic or impersonate a legitimate speaker, leading to false acceptances. Spoofing attacks involve presenting artificial or manipulated speech signals to deceive the verification system. Examples of spoofing attacks include playback attacks, voice conversion, and speech synthesis. Detecting and mitigating impostor and spoofing attacks requires the integration of anti-spoofing techniques, such as analyzing speech characteristics, using multi-modal information (e.g., facial or behavioral cues), or employing specialized detectors to distinguish genuine from manipulated

speech signals.

6.5 Computational Complexity:

Speaker verification systems often involve complex processing tasks, including feature extraction, modeling, and decision-making, which can be computationally demanding. Real-time or near real-time performance is crucial for many speaker verification applications, such as access control or authentication systems. Efficient algorithms, optimization techniques, and hardware acceleration methods are needed to address the computational complexity challenge and ensure efficient and scalable speaker verification solutions.

Addressing these challenges in a speaker verification project requires careful consideration of the system design, choice of algorithms and models, dataset preparation, and evaluation methodologies. Innovative solutions and ongoing research efforts are aimed at improving the robustness, reliability, and efficiency of speaker verification systems in real-world applications.

VII Performance Evaluation and Datasets:

7.1 Evaluation Metrics:

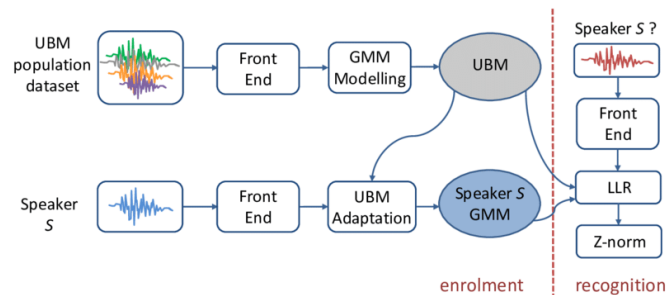
In speaker verification, several evaluation metrics are used to assess the performance of the system. The choice of evaluation metrics depends on the specific requirements and goals of the project. Some commonly used evaluation metrics in speaker verification include:

Equal Error Rate (EER): The point at which the false acceptance rate (FAR) and false rejection rate (FRR) are equal. EER is a single, threshold-independent metric of system performance. **Detection Cost Function (DCF):** A metric that combines the errors made by the system with the associated costs of these errors. DCF considers both the false acceptance and false rejection errors and allows for the incorporation of different costs for each error type.

Precision and Recall: These metrics are commonly used in binary classification problems. Precision represents the proportion of correctly identified positive instances (genuine speakers) out of the total instances identified as positive. Recall measures the proportion of correctly identified positive instances out of the actual positive instances.

Receiver Operating Characteristic (ROC) curve: A graphical illustration of the trade-off between the true positive rate (TPR) and the false positive rate (FPR) at various judgement thresholds. The area under the ROC curve (AUC) is a common assessment statistic, where a higher AUC indicates better system performance.

Detection Error Tradeoff (DET) curve: A plot of the false acceptance rate (FAR) against the false rejection rate (FRR) on logarithmic scales. The DET curve provides a visualization of system performance and allows for easy comparison across different systems.



7.2 Commonly Used Datasets:

Several publicly available datasets are commonly used for training, development, and evaluation of speaker verification systems. Some well-known datasets include:

NIST SRE: The National Institute of Standards and Technology Speaker Recognition Evaluation (NIST SRE) datasets are widely used in speaker verification research. These datasets contain a large collection of speech samples from multiple speakers recorded under various conditions.

VoxCeleb: The VoxCeleb dataset consists of a vast amount of speech data collected from celebrities in different languages. It provides a diverse set of speakers and variations in recording conditions, making it suitable for training and evaluating speaker verification models.

LibriSpeech: The LibriSpeech dataset is a collection of read speech from audiobooks, containing speech samples from a large number of speakers. It covers a wide range of speech variations and is often used for training and evaluating speaker verification models.

TIMIT: The TIMIT dataset is a widely used resource for speech research, including speaker verification. It contains speech samples from multiple speakers of different dialects, ages, and genders.

RedDots: The RedDots dataset is a collection of multi-modal biometric data, including speech, video, and other biometric modalities. It is designed for research in multi-modal speaker verification and anti-spoofing.

7.3 Benchmark Results:

Benchmark results for speaker verification systems are typically reported in terms of the evaluation metrics mentioned earlier, such as EER, DCF, or accuracy. These results provide a comparative analysis of different systems or approaches on a specific dataset or evaluation protocol.

Benchmark results can be found in research papers, technical reports, or challenge evaluations in the field of speaker verification. For example, the NIST Speaker Recognition Evaluation (SRE) reports provide detailed benchmark results of different systems on the NIST SRE datasets. Similarly, challenges like the Speaker Recognition Evaluation (SRE) organized by the International Speech Communication Association (ISCA) provide benchmark results and facilitate the comparison of different systems and algorithms.

VIII CONCLUSION:

8.1 Summary of Findings:

In this paper, we have explored the field of speaker verification, highlighting various aspects related to feature extraction, speaker modeling approaches, decision-making techniques, challenges, performance evaluation, and datasets.

We discussed the motivation behind speaker verification, its purpose in access control is one example of an application and authentication, and the key components involved in the process. We examined different feature extraction techniques commonly used in speaker verification, including Mel-Frequency Cepstral Coefficients (MFCCs), Linear Predictive Coding (LPC), Perceptual Linear Prediction (PLP), and Deep Learning-based features. These techniques play a crucial role in capturing speaker-specific characteristics and transforming speech signals into discriminative representations. Furthermore, we explored various speaker modeling approaches, including Gaussian Mixture Models (GMMs), Support Vector Machines (SVMs), Deep Neural Networks (DNNs), Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Speaker Embeddings. These modeling techniques enable the development of robust and accurate speaker verification systems by capturing and modeling speaker-specific information. We also discussed decision-making techniques such as score normalization, threshold selection, Bayesian approaches, and fusion of multiple models. These techniques contribute to enhancing the reliability and efficiency of speaker verification systems. Moreover, we addressed the challenges in speaker verification, including data availability and collection, variability in speech signals, channel and environmental conditions, impostor detection, and spoofing attacks, as well as computational complexity. Understanding and mitigating these challenges are essential for building effective and robust speaker verification systems. In addition, we examined the evaluation metrics commonly used to evaluate the effectiveness of speaker verification systems, the commonly used datasets for training and evaluation purposes, and the importance of benchmark results in comparing different approaches and algorithms.

8.2 Key Takeaways:

From our exploration, the following key takeaways can be summarized: Speaker verification is an important technology that enables secure and reliable access control and authentication systems. Feature extraction techniques such as MFCCs, LPC, PLP, and deep learning-based features are instrumental in capturing speaker-specific information from speech signals. Speaker modeling approaches such as GMMs, SVMs, DNNs, CNNs, RNNs, and speaker embeddings enable the development of accurate and robust speaker verification systems.

Decision-making techniques such as score normalization, threshold selection, Bayesian approaches, and fusion of multiple models contribute to enhancing the decision-making process and system performance. Speaker verification faces challenges such as data availability, variability in speech signals, channel and environmental conditions, impostor detection, and computational complexity. Addressing these challenges is crucial for building effective systems. Evaluation metrics, commonly used datasets, and benchmark results provide means to assess and compare the performance of different speaker verification systems.

8.3 Importance of Speaker Verification:

Speaker verification plays a vital role in various applications

where secure and reliable authentication is required. It offers a non-intrusive and convenient method for verifying individuals based on their unique voice characteristics. Speaker verification systems can enhance security measures in access control systems, financial transactions, telephone banking, voice assistants, and forensic investigations.

REFERENCES

- [1] 1.Reynolds, D., & Campbell, W. M. (2008). Speaker verification using adapted Gaussian mixture models. *Digital Signal Processing*, 18(3), 275-284.
- [2] 2.Kenny, P., Ouellet, P., & Dehak, N. (2010). A study of interspeaker variability in speaker verification. *IEEE Transactions on Audio, Speech, and Language Processing*, 18(4), 868-878.
- [3] 3.Doddington, G., Przybocki, M., Martin, A., & Reynolds, D. (2000). The DET curve in assessment of detection task performance. *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 3, 1385-1388.
- [4] 4.Sell, G., McCree, A., & Sheikhzadeh, H. (2018). Deep learning for speaker recognition: An overview. *IEEE Signal Processing Magazine*, 35(2), 81-102.
- [5] 5.Snyder, D., Garcia-Romero, D., & Povey, D. (2018). Deep neural network embeddings for text-independent speaker verification. *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 4874-4878.
- [6] 6.Brümmer, N., & Burget, L. (2011). Joint factor analysis versus eigenchannels in speaker recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(4), 788-798.
- [7] 7.Srivastava, R., Grezl, F., & Cernocký, J. (2019). Convolutional neural networks for speaker recognition. *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 5906-5910.
- [8] 8.Reynolds, D., Quatieri, T., & Dunn, R. (2000). Speaker verification using adapted Gaussian mixture models. *Digital Signal Processing*, 10(1-3), 19-41.
- [9] 9.Zeinali, H., & Hansen, J. H. (2016). Deep recurrent neural network-based feature mapping for speech recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 24(9), 1572-1583.
- [10] 10.Doddington, G., Liggett, W., Martin, A., Przybocki, M., & Reynolds, D. (1997). Sheep, goats, lambs and wolves: A statistical analysis of speaker performance in the NIST 1996 speaker recognition evaluation. *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 1, 357-360.
- [11] 11.An Effective Dereverberation Algorithm by Fusing MVDR and MCLP Open Access Tan, F., Bao, C., Liu, T.2022 IEEE International Conference on Signal Processing, Communications and Computing, ICSPCC 2022
- [12] 12. DA-VAD: UNPAIRED ADVERSARIAL DOMAIN ADAPTATION FOR NOISE-ROBUST VOICE ACTIVITY DETECTION Kim, T., Chang, J., Ko, J.H. 2022 ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing – Proceedings
- [13] 13. Audiovisual speaker indexing for Web-TV automations Vryzas, N., Vrysis, L., Dimoulas, C.2021 Expert Systems with Applications
- [14] 14. A novel voice activity detection algorithm using modified global thresholding Elton, R.J., Mohanalin, J., Vasuki, P. 2021 International Journal of Speech Technology
- [15] 15. Voice activity detection using generalized exponential kernels for time and frequency domains
- [16] 16 Soares, A.D.S.P., Parreira, W.D., Souza, E.G., Nascimento, C.D.D., Almeida, S.J.M.D.2019 IEEE Transactions on Circuits and Systems I: Regular Papers