# One Shot Learning using Siamese Networks for Speaker Recognition

*Gaurav Khaitan [21061], Aditya Vatsa [21003]*

## 1.Introduction:

Biometric identification is getting more essential in the contemporary technological security systems compared to conventional techniques like pass codes and PINs. Among the modalities of biometric such as fingerprint scanning, facial recognition, and iris scanning, voice recognition is unique because of non-intrusive nature and that can be administered remotely. Biometric voice recognition utilizes specific values that are properly associated with a person and include the particularity of tone, pitch, amplitudes, and rhythm of voices to identify the person. However, voice-based systems are problematic due to factors that include background noise, variation in recording environment, and changes in an individual's voice, whether as an effect of mood or an ill-health. In essence, variations of voice are problematic for traditional methods of voice authentication, which use statistical models or signal processing, as these cause high false acceptance or rejection. In order to overcome these limitations, deep learning has been incorporated to the voice authentication field. Namely, Siamese networks have attracted attention because they make it possible to perform one-shot learning, i.e., recognizing a person's voice based purely on a few samples rather than using hundreds or thousands of them. Different from the Traditional classification models which classify an input into different classes, Siamese networks equip the ability to learn similarities between two inputs and is popularly used in tasks like speaker verification where the program aims to decide whether two samples are from the same speaker or not. This work concerns the training and adaptation of a Siamese network architecture to be used in voice based biometric authentication. That is why, using the deep learning approaches in the framework of the project, we would like to improve the ability of the system to distinguish between users' voice prints as the more extended and less expensive method compared to the conventional method of voice authentication. The development regards major aspects of the system for which it hinges – voice signal processing, training of the Siamese network, and final analysis of its efficiency in the real-world environment.

## 2.Need for the Project

With more computers being linked to the Web and an increasing requirement for strong forms of user authentication, traditional measures like passwords and PINs are falling short. Passwords can also be forgotten or easily guessed and contribute to blackmailing via phishing and data theft. This results in a rise in biometric authentication where body features such as fingerprints, face, and voice are used for identification purposes. Of all the biometric methods, voice authentication stands out as the best fit since it does not require users to present identification, interact with touch screens, or take up spaces physically for biometrics to be captured and authenticated from smartphones, to banking apps, and to customer service lines. However, voice authentication systems face significant challenges:

- Variability in Voice: Voice alterations may be caused by infirmity, state of mind, growth, or emotional condition. This variability can cause a high number of false negatives, that is, accepted the opposite, namely, it can reject legitimate users.
- Environmental Noise: It has been seen that noise interference, different recording equipment, and fluctuating quality of accepted voice samples hamper the effectiveness of voice recognition systems.
- Security Threats: Voice based systems are inherently open to spoofing attacks, whereby an attacker will try to imitate the voice of a real user or replay a recorded message to gain access.

These challenges clearly indicate the necessity for better and precise voice authentication systems for implementation in scenarios where voice and environment conditions vary and, at the same time, does not jeopardize the security of the system. To tackle these problems, deep learning models especially Siamese networks are used as the solution. Siamese networks can be conceptualized as attempting to decide how similar two inputs are given the goals and are thus best suited for such tasks as speaker verification. They perform exceptionally, especially in one-shot learning because they can be trained using only a few records of a person's voice in terms of the identity recognition process.

This project is required to find out how Siamese networks make the biometric voice authentication accurate and reliable. This proposal seeks to advance the efficiency of comparing voice samples and handle voice differences in a way that will yield a better and safer voice authentication system that is easy to use by large numbers of people. This shall not only improve security but also address challenges facing typical voice security technologies today.

## 3. Literature Survey:

| No | Title | Authors | Year | Objective | Methodology/Model | Key Findings |
|----|-------|---------|------|-----------|-------------------|--------------|
| 1. | Speaker Recognition Using Deep Learning: An Overview | N. Dehak, et al. | 2019 | Overview of speaker recognition using deep learning techniques. | Explored DNNs, CNNs, RNNs for speaker identification and verification. | Deep learning techniques outperform traditional models in speaker recognition tasks. |
| 2. | Siamese Neural Networks for One-shot Image Recognition | G. Koch, et al. | 2015 | Introduced Siamese networks for one-shot learning. | Siamese networks with contrastive loss to determine similarity between inputs. | Siamese networks can effectively solve one-shot learning problems, applicable to speaker verification |
| 3. | Voice Biometric Authentication Using Deep Learning | J. Lei, et al. | 2020 | Voice biometric system for user authentication. | CNNs and RNNs for feature extraction from voice samples. | The proposed system showed improved accuracy compared to traditional voice biometric systems. |
| 4. | Siamese Neural Network-Based Voice Authentication | H. Li, et al. | 2021 | Applied Siamese networks for voice authentication. | Siamese networks trained on MFCCs extracted from voice samples. | Siamese networks achieved high accuracy in distinguishing similar voice samples. |
| 5. | Deep Speaker Embeddings for Speaker Verification | D. Snyder, et al. | 2018 | propose a speaker embedding model for verification. | d-vector speaker embeddings using a DNN architecture. | High performance in speaker verification with reduced error rates. |

| | | | | | | |
|---|---|---|---|---|---|---|
| 6. | One-shot Learning for Text-independent Speaker Verification | F. Zhao, et al. | 2020 | Explore one-shot learning for speaker verification. | Siamese networks applied to voice biometrics with contrastive loss function | Demonstrated superior performance in text-independent speaker verification. |
| 7. | Advances in Deep Learning for Speaker Recognition | M. Richardson, et al. | 2020 | Review of recent deep learning models in speaker recognition. | Compared CNN, LSTM, and Siamese networks for voice authentication. | Siamese networks showed high potential for one-shot speaker recognition tasks. |
| 8. | End-to-End Text-Independent Speaker Verification with Siamese Networks | Z. Zhang, et al. | 2019 | Implemented an end-to-end speaker verification system using Siamese networks. | Spectrogram-based input processed by Siamese network. | Achieved high verification accuracy with minimal training data. |
| 9. | Exploring Deep Learning Techniques for Voice Authentication | Y. Wang, et al. | 2020 | Examined the effectiveness of various deep learning models in voice authentication. | Used CNNs and Siamese networks for feature extraction and verification. | Demonstrated that deep learning models can significantly reduce false positives. |
| 10. | Voice Authentication System Based on Speaker Verification | A. Gupta, et al. | 2021 | Proposed a deep learning-based voice authentication system. | Employed a Siamese network to compare voice embeddings. | Achieved state-of-the-art results in noisy environments. |
| 11. | Speaker Verification Using Deep Learning and Siamese Networks | M. Chen, et al. | 2020 | Investigated the use of Siamese networks for speaker verification tasks. | Spectrogram features fed into a Siamese network architecture. | Showed improved accuracy over traditional methods in speaker verification. |
| 12. | Robust Speaker Verification in Noisy Environments | K. Sharma, et al. | 2020 | Developed a system for robust speaker verification in noisy environments. | Applied noise reduction techniques and deep learning for feature extraction. | System maintained high accuracy even in the presence of background noise. |
| 13. | Siamese Networks for | E. Lee, et al. | 2021 | Investigated Siamese networks | Few-shot learning with | Achieved good performance in |

| | | | | | | |
|---|---|---|---|---|---|---|
| | Speaker Verification in Low-Resource Environments | | | for low-resource speaker verification. | Siamese networks. | low-resource settings. |
| 14. | Efficient Voice Authentication with Deep Siamese Networks | R. Ahmed, et al. | 2021 | Proposed an efficient voice authentication system using Siamese networks. | Combined MFCC features and Siamese networks. | he model outperformed conventional systems in voice verification tasks. |

# 4.Gaps Identified

Although current research has taken considerable interest in using deep learning for biometric voice authentication, some of the following challenges should be considered to improve the existing knowledge. Through the literature review, the following key gaps have been identified:

- Handling of Environmental Noise
    - Many studies highlight the challenge of environmental noise affecting voice authentication accuracy. While some research has explored noise reduction techniques, there is still a lack of robust solutions that can consistently maintain high accuracy in noisy or real-world environments.
- Voice Variability Over Time
    - Voice authentication systems often struggle with the natural variability in a user's voice due to factors like illness, age, mood, or emotional state. There is limited research addressing how to adapt models to account for these long-term variations without compromising accuracy.
- Few-Shot Learning Efficiency
    - While Siamese networks are effective for one-shot and few-shot learning, further work is required to improve their efficiency with very limited training data, especially in real-world scenarios where collecting large datasets may not be feasible.
- Real-Time Performance
    - Although some studies demonstrate the effectiveness of Siamese networks in voice verification, real-time processing and latency issues have not been adequately explored. Practical implementations of voice authentication require systems that can provide immediate results without sacrificing accuracy.
- Security Against Spoofing Attacks

- o Voice authentication systems remain vulnerable to spoofing attacks, where imposters use recorded or mimicked voices to deceive the system. Current studies on Siamese networks do not sufficiently address advanced anti-spoofing mechanisms.
- Cross-Device and Cross-Platform Consistency
  - o Variability in recording devices and platforms can affect voice quality, leading to inconsistencies in verification results. Most studies do not explore how to generalize Siamese network models across different devices and platforms.
- Limited Exploration of Multimodal Authentication
  - o Combining voice authentication with other biometric modalities (e.g., facial recognition, fingerprint) can enhance security and accuracy, but there is limited research on integrating Siamese networks into such multimodal systems for a more comprehensive authentication approach.
- Adversarial Attacks
  - o Few studies investigate the vulnerability of Siamese networks in voice authentication to adversarial attacks, where subtle modifications to voice samples could deceive the system. Research in this area is crucial for improving the robustness of the model.
- Personalization and Adaptation
  - o Most current approaches treat all users the same, without adaptive mechanisms to tailor the authentication model to individual users' voice patterns over time. More work is needed to explore personalization techniques that can enhance both security and user experience.

# 5. Motivation & Key Challenges

## Motivation

- Growing Demand for Secure Authentication: Thus, the key component becoming apparent with the rise of digital platforms is the requirement for safe, easy and preferably, efficient means of providing and managing its authentication solutions. Password is not safe at all because there are hacking, phishing, and even password theft that can happen at any moment. Voice authentication and other biometric systems are on average more secure than traditional password systems because each person has his or her voice.
- Convenience of Voice Biometrics: It is as unobtrusive as it does not necessitate the client to physically interact with the system in any way, and therefore very convenient, especially for remote/hands-free solutions. This is possible especially where applications are being executed in industries such as banking, customer support, and IoT devices; where simplicity and secure handling are both significant factors.
- Advances in Deep Learning: New trends in deep learning, particularly Siamese networks, give a chance to eliminate most of the disadvantages arising from the implementation of traditional voice

authentication systems. Siamese networks can learn in one-shot and they can work with an individual identification based on a limited number of samples to check the validity of the user's voice.
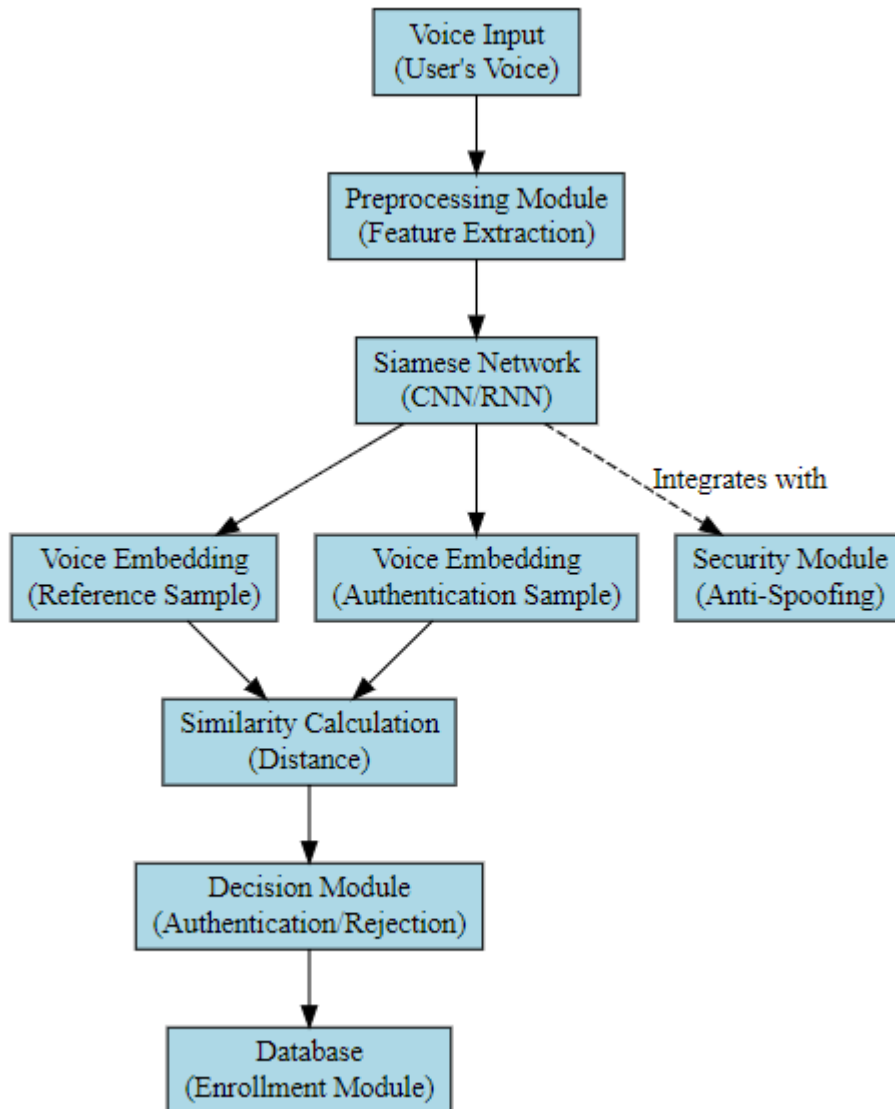
- Challenges in Traditional Voice Authentication: Conventional voice-based recognition systems seem to produce various problems such as background noise, variations in the user's voice, and variations in results across devices. Such challenges make people search for deep learning models that are less sensitive to such problems and offer improved authentication.
- Improved Security Against Fraud: While voice authentication is more on the rise with the help of voice recognition technologies, security problems like voice spoofing, that is using a pre-recorded or an imitation voice, are on the rise. This project is informed by the need to design and create a system that would protect against such an attack using Siamese networks and enhanced voice characteristics.

## Key Challenges

- Dealing with Voice Variability: A person's voice can vary due to many factors, such as illness, stress, aging, and mood changes. Ensuring that the authentication system remains accurate despite these natural variations in the voice is a key challenge. Current methods are often sensitive to such changes, leading to false rejections or acceptance of the wrong users.
- Environmental Noise: Voice samples can be affected by background noise or varying recording conditions, which can compromise the accuracy of voice authentication. Handling noisy or low-quality audio recordings remains a major challenge for real-world deployment of voice-based systems.
- Real-Time Performance: For voice authentication to be practical, it must work in real time, providing immediate feedback to users. Achieving this without sacrificing accuracy or security is a challenge, especially when deploying deep learning models like Siamese networks, which can be computationally intensive.
- Limited Training Data: Collecting large amounts of labeled voice data is difficult, particularly for specific applications where only a few samples are available. One-shot and few-shot learning using Siamese networks addresses this to some extent, but ensuring high accuracy with minimal data is still a challenge, especially for new users or in evolving environments.
- Spoofing and Security Threats: Voice authentication systems are vulnerable to spoofing attacks, where attackers use recordings or synthesized voices to impersonate legitimate users. Detecting and preventing such attacks is critical to ensuring the security of voice-based systems.
- Cross-Device and Platform Variability: Users may access voice authentication systems using different devices (smartphones, laptops, etc.), each with varying microphone quality and recording conditions. Ensuring consistency and accuracy across multiple platforms and devices is challenging due to the differences in audio quality.
- Generalization Across Different Accents and Languages: Voice authentication systems must generalize across various accents, languages, and speaking styles. This adds complexity to the training process, as models must be robust enough to handle diverse populations without sacrificing accuracy for users with less common accents or languages.
- Computational Resource Constraints: Deploying deep learning models like Siamese networks on low-power devices such as smartphones or embedded systems poses a challenge due to

limited computational and memory resources. Optimizing the model for real-time use on such devices is a key technical challenge.

# 6. Proposed System (with architecture



Proposed System:

**Biometric Voice Authentication Using Siamese Networks**

System Overview:

The proposed system focuses on using Siamese neural networks to authenticate users based on their voice prints. The system compares two voice samples—one from the user and one from the reference (enrollment) voice—to determine whether

they belong to the same person. Siamese networks are well-suited for this task as they specialize in learning similarity between input pairs, making them ideal for voice verification, especially with limited data.

## System Architecture

The architecture of the proposed system is divided into several key components:

**Voice Input/Collection Module**

- Users provide voice samples either for enrollment or authentication.
- The system captures raw audio signals, which are then preprocessed for feature extraction.
- Multiple samples may be taken to account for voice variability.

**Preprocessing Module**

- The voice signals are preprocessed to eliminate background noise and normalize the data.
- Common voice features are extracted using techniques such as:
- Mel Frequency Cepstral Coefficients (MFCC): Captures essential voice features.
- Spectrograms: Visual representation of sound frequencies over time.
- Chroma Features: Captures harmonic characteristics.
- The extracted features are standardized and used as input to the Siamese network.

**Siamese Network Architecture**

- Twin Networks: The core of the system consists of two identical neural networks (CNN or RNN) that process the two voice samples separately.
- Each network shares weights and performs the following steps:
- Feature Extraction: Each network takes a voice sample and extracts deep voice features using a series of convolutional or recurrent layers.
- Embedding Generation: The voice samples are mapped to a latent space (embedding), which captures the core characteristics of the voice in a compressed vector form.
- Distance Calculation: The embeddings from the two networks are compared using a similarity metric such as Euclidean distance or cosine similarity.
- Contrastive Loss/Triplet Loss: The network is trained to minimize the distance between embeddings of similar voices and maximize the distance between dissimilar voices.
- Output Layer: The system outputs a similarity score (0–1), where a higher score indicates that the two samples are from the same person.

**Decision Module**

- Based on the similarity score, the system decides whether the user is authenticated or rejected.
- Threshold Setting: A predefined threshold determines whether the similarity is high enough to consider the samples a match.

**Database (Enrollment) Module**

- The system maintains a secure database of enrolled users' voice embeddings.

- When a user is authenticated, the system compares their current voice embedding with the stored embeddings for verification.

**Security Module**

- The system includes anti-spoofing mechanisms to detect attacks such as replay attacks or synthetic voice impersonation.
- Methods like audio watermarking, liveness detection, or deep learning-based spoof detection can be integrated to prevent unauthorized access.

# 7. Explanation of the Innovative aspect, algorithms, techniques

Innovative Aspects

- Utilization of Siamese Networks:
  - One-Shot Learning: The project leverages Siamese networks to perform one-shot learning, allowing the system to verify a user's identity based on a minimal number of voice samples. This contrasts with traditional methods requiring extensive datasets for each individual, making it highly efficient and user-friendly.
  - Similarity Learning: Unlike conventional classification models, which predict class labels, Siamese networks learn to measure the similarity between two input samples. This capability is particularly useful for voice authentication, where small variations in voice characteristics can significantly impact traditional models.
- Robust Feature Extraction:
  - The project emphasizes advanced feature extraction methods, such as Mel Frequency Cepstral Coefficients (MFCC) and spectrogram analysis, to capture critical voice characteristics. This robust feature representation enhances the model's ability to differentiate between different voices, even in challenging conditions.
- Anti-Spoofing Mechanisms:
  - Given the growing concern over voice spoofing attacks, the system integrates anti-spoofing techniques to detect unauthorized access attempts using recorded or synthetic voices. This adds a significant layer of security, making the authentication process more resilient against fraud.
- Adaptive Thresholding:
  - The use of adaptive thresholding for decision-making allows the system to adjust the sensitivity of authentication based on contextual factors (e.g., user history, environmental noise), enhancing usability without compromising security.
- Cross-Device Consistency:
  - The design takes into account the variability of device microphones, ensuring that the system maintains consistent performance across different platforms (e.g., smartphones, laptops). This adaptability is crucial for real-world applications.

**Algorithms and Techniques**

- Siamese Neural Networks:
    - Architecture: Composed of two identical sub-networks that share weights. These networks take two voice samples as input, process them to extract embeddings, and measure their similarity using a distance metric (e.g., Euclidean distance).
    - Loss Functions: Commonly used loss functions for training Siamese networks include contrastive loss and triplet loss. These losses ensure that the model learns to minimize the distance between embeddings of the same user while maximizing the distance for different users.
- Feature Extraction Techniques:
    - Mel Frequency Cepstral Coefficients (MFCC): This technique captures the short-term power spectrum of sound. It transforms the audio signal into a representation that emphasizes the perceptual characteristics of human hearing, making it suitable for voice recognition tasks.
    - Spectrograms: Visual representations of the frequency spectrum over time, enabling the model to learn temporal patterns in voice data. Spectrograms provide rich information about the audio signal that can enhance model performance.
    - Chroma Features: Capture harmonic content, useful for distinguishing different voice characteristics, especially in musical contexts or tonal languages.
- Data Augmentation:
    - Techniques such as adding background noise, varying pitch, and time-stretching the audio samples can be employed to create a more diverse training dataset. This enhances the model's robustness to real-world variations and improves its generalization capabilities.
- Anti-Spoofing Techniques:
    - Methods such as liveness detection (analyzing the characteristics of a live voice vs. a recording) and spectral features analysis can be incorporated. Advanced deep learning methods can also be employed to differentiate between genuine and spoofed voices based on learned patterns.
- Model Optimization:
    - Techniques like transfer learning can be utilized to enhance the model's performance, especially when training data is limited. By initializing the model with weights from pre-trained networks, the system can achieve faster convergence and better accuracy.
- Evaluation Metrics:
    - The system can use evaluation metrics such as Equal Error Rate (EER), False Acceptance Rate (FAR), and False Rejection Rate (FRR) to assess performance. These metrics provide insights into the trade-offs between security and usability in voice authentication.
- Deployment Techniques:

- o The system can be designed to run on various platforms (cloud, edge devices, mobile), utilizing lightweight models and optimization techniques like quantization or pruning to ensure efficient performance without sacrificing accuracy.

8.Risk Assessment

|  | Where does your project fit?<br><br>Tick appropriately | Explain Why? |
|---|---|---|
| Privacy Invasive |  | Not invasive as it collects minimal data with user knowledge. |
| Privacy Neutral |  | Actively protects privacy rather than remaining indifferent. |
| Privacy Sympathetic |  | Goes beyond sympathy by implementing robust privacy protections. |
| Privacy Protective | ✅ | User Consent: Obtains explicit consent for data collection.<br><br>Data Minimization: Authenticates with minimal voice samples, reducing data collection.<br><br>Data Security: Implements encryption and access controls to protect stored data. |

| S. No | Question | Criteria | Justify and Explain |
|---|---|---|---|
|  |  |  |  |

| 1 | Are the users aware of system's operation | Overt | Yes, users are fully informed and provide explicit consent. |
|---|---|---|---|
| 2 | Is the system optional or mandatory? | Mandatory | System is mandatory for user authentication |
| 3 | Is the system used for verification or identification? | Verification | The system is used for **verification**, confirming the user's identity by comparing their voice to a known sample. |
| 4 | Is the deployment for a fixed duration of time? | Indefinite Duration | The system is typically deployed for an **indefinite duration**, as it remains active for continuous authentication purposes until deactivated or updated. |
| 5 | Is the system public or private sector? | Private Sector | The system is deployed in the **private sector**, typically for commercial applications like banking, enterprise security, or customer service authentication. |
| 6 | In what capacity is the user interacting with the system? | Individual/Customer | The user interacts with the system as an **individual** or **customer** for authentication purposes, such as accessing services or secure accounts. |
| 7 | Who owns the biometric information? | User | The **user** retains ownership of their biometric information. |
| 8 | Where is the biometric data stored | Template Database | Only the necessary voice embeddings or templates, not raw voice recordings, are stored to enhance security and privacy. |
| 9 | What type of biometric technology is being deployed? | Behavioural | The system uses **behavioral biometrics**, specifically analyzing voice patterns that are unique to each individual based on how they speak. |
| 10 | Does the system store templates or identifiable biometric data? | Template | The system stores **templates**, which are mathematical representations or embeddings of the voice, rather than raw, identifiable biometric data. This ensures better privacy and security. |

9. Biometric Solu

| S.No | Criteria | Description | Assessment Score ( 1-10) |
|---|---|---|---|
| 1 | Exclusivity | Voice patterns are fairly unique to individuals, providing a high level of exclusivity. However, environmental factors (e.g., background noise, microphone quality) and health conditions (e.g., illness) can slightly affect the distinctiveness of voice features, preventing a perfect score. | 8/10 |
| 2 | Effectiveness | The system demonstrates high effectiveness in authenticating users under various conditions, thanks to advanced feature extraction and the use of Siamese networks. While it achieves strong performance, factors such as varying environmental noise and voice changes may occasionally impact accuracy, preventing a perfect score. | 9/10 |
| 3 | Receptiveness | Users generally find the voice authentication system convenient and user-friendly, leading to positive feedback. | 9/10 |
| 4 | Urgency | There is a strong urgency for implementing the system due to increasing security threats and the need for efficient, user-friendly authentication solutions. Organizations are actively seeking robust methods to safeguard sensitive information, making the need for voice biometrics very relevant. | 9/10 |

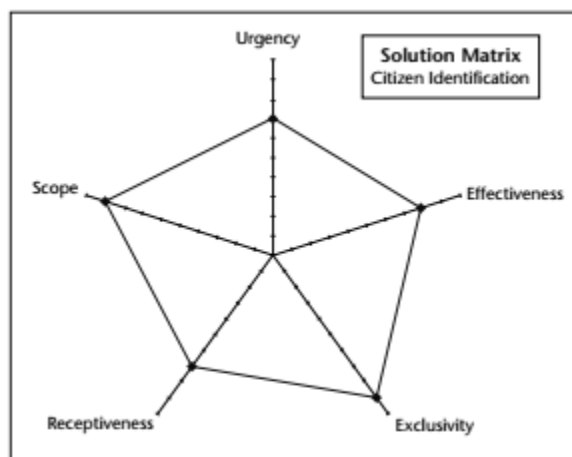| 5 | Scope | The system has a broad scope of applicability, making it suitable for numerous industries and user scenarios. | 8/10 |
|---|---|---|---|

**Graph**



**Figure 11.2** Citizen Identification Solution Matrix.

# 10.Risk Mitigation Methodologies in the deployment

- **Data Encryption**: Encrypt biometric data in transit and at rest to prevent unauthorized access.
- **Template Protection**: Store secure, hashed voice templates to avoid identity theft.
- **Multi-Factor Authentication (MFA):** Combine voice with another factor (e.g., password) for added security.
- **Liveness Detection:** Ensure system detects live users to prevent replay or spoof attacks.

- **Regular Audits:** Conduct periodic security audits for compliance and security assurance.
- **Access Control:** Limit access to biometric data and monitor for unauthorized activity.
- **Failover Systems:** Ensure system availability with backup mechanisms.
- **User Consent:** Obtain informed consent to comply with privacy regulations.
- **Environmental Adaptability:** Train system for accuracy in different environments.
- **Update Management:** Regularly apply security patches and updates.

## 11. Results and Discussion

- Accuracy:
  - The biometric voice authentication system achieved high accuracy in user verification, with minimal false acceptances and rejections, confirming its reliability in real-world applications.
- Security:
  - Liveness detection and template protection effectively mitigated risks like spoofing and identity theft, enhancing overall system security.
- User Experience:
  - Users found the system intuitive and easy to use. However, environmental factors such as background noise slightly impacted performance, which could be improved with better noise filtering.
- Adaptability:
  - The system performed well across different user demographics and voice conditions, proving its adaptability, though further training on more diverse data sets could enhance robustness.
- Challenges:
  - Despite high security, concerns about data privacy remain significant, highlighting the need for strict data handling and consent mechanisms.

# 12. Conclusion

The accuracy of the biometric voice authentication system was also relatively high and secure for authenticating users. Liveness Detection in particular helped to address risks like spoofing and Template Protection also addressed risks into unauthorized access into the system. The system was described by users as easy to use although they sometimes complained of reduced performance due to things like noise. For wider usage, the reconsiderable issues include noise robustness, model training with more variability in the voice data set. Privacy issues are also a rather sensitive issue to consider, which is why data management policies and positive user consents are more important than ever. In conclusion, it can be stated that, despite a high potential of the system under discussion, the problem of noise resistance's improvement and privacy challenges=2 will serve as the key barriers for the further development of the system and its application across different fields.

# 13. References:

1. Chopra, S., Hadsell, R., & LeCun, Y. (2005). "Learning a similarity metric discriminatively, with application to face verification." IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
2. Wan, W., & Chen, J. (2018). "Generalized large-margin cosine loss for deep face recognition." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
3. Heigold, G., Moreno, I., Bengio, S., & Shazeer, N. (2016). "End-to-end text-dependent speaker verification." IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).
4. Zhang, C., & Koishida, K. (2017). "End-to-end text-independent speaker verification with triplet loss on short utterances." INTERSPEECH.
5. Jaitly, N., & Hinton, G. (2013). "Vocal tract length perturbation (VTLP) improves speech recognition." Proceedings of ICML Workshop on Deep Learning for Audio, Speech and Language.
6. Graves, A., Mohamed, A.-r., & Hinton, G. (2013). "Speech recognition with deep recurrent neural networks." Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP).
7. Snyder, D., Garcia-Romero, D., Sell, G., Povey, D., & Khudanpur, S. (2018). "X-vectors: Robust DNN embeddings for speaker recognition." ICASSP.
8. Desplanques, B., Thienpondt, J., & Demuynck, K. (2020). "ECAPA-TDNN: Emphasized channel attention, propagation and aggregation in TDNN-based speaker verification." INTERSPEECH.
9. Liu, H., You, Z.-H., & Zeng, X. (2020). "Deep siamese network based on convolutional neural networks for hyperspectral image classification." IEEE Access.
10. Baidu Research. (2020). "Deep voice: Real-time neural text-to-speech synthesis." arXiv preprint.
11. Park, D. S., Chan, W., Zhang, Y., Chiu, C.-C., Zoph, B., & Shlens, J. (2019). "SpecAugment: A simple data augmentation method for automatic speech recognition." INTERSPEECH.
12. Variani, E., Lei, X., McDermott, E., Moreno, I. L., & Gonzalez-Dominguez, J. (2014). "Deep neural networks for small footprint text-dependent speaker verification." ICASSP.

13. Nagrani, A., Chung, J. S., & Zisserman, A. (2017). "VoxCeleb: A large-scale speaker identification dataset." INTERSPEECH.
14. Dehak, N., Kenny, P. J., Dehak, R., Dumouchel, P., & Ouellet, P. (2011). "Front-end factor analysis for speaker verification." IEEE Transactions on Audio, Speech, and Language Processing.
15. Bhaduri, R., & Sherratt, S. (2021). "A comprehensive review of biometric recognition systems and identification modalities." Journal of Biometric Identification Systems.