# CRNN Algorithm and MFCC Feature Extraction in Classifying *Hijaiyah* Letter Pronunciation: A Systematic Literature Review

Mohammad Putra Fauzan Fatah
*Department of Informatics*
UIN Sunan Gunung Djati, Indonesia
Putrafauzan32@gmail.com

*Abstract*— The ability to read the *hijaiyah* letters correctly is an important foundation in learning the Qur'an. However, the low Qur'an literacy in Indonesia indicates the need for technological innovation to support the learning process. This article presents a systematic literature review on the application of the Convolutional Recurrent Neural Network (CRNN) algorithm and the Mel-Frequency Cepstral Coefficients (MFCC) feature extraction technique in speech classification, with a focus on their potential implementation for recognizing the pronunciation of the *hijaiyah* letters. The analysis was conducted based on ten relevant studies selected using the PRISMA method. The results of the study indicate that MFCC is effective in representing the phonetic characteristics of sounds, including in the Arabic context. Meanwhile, CRNN has proven superior in managing audio data with tempo and sequential structure. The combination of the two has strong potential to build an accurate and adaptive *hijaiyah* letter sound classification system, especially in supporting speech-based tajweed learning. This study provides a conceptual basis for the development of an artificial intelligence-based Qur'an learning application.

*Keywords- Audio Classification, Al-Qur'an literacy, CRNN, Hijaiyah Letters, MFCC, Voice Recognition.*

## I. INTRODUCTION

The ability to read the Quran properly and correctly is one of the fundamental competencies in the life of a Muslim. Learning the Quran also means understanding how to pronounce each letter, training the tongue to pronounce the letters according to their place of origin, combining letters with each other, reading with a variety of long and short tones, reading techniques by omitting the sound of a letter and connecting it with the next letter, as well as understanding the heaviness and lightness when pronouncing letters, producing hissing sounds or not, and learning the stop signs in the sentences read.

According to Tombak Alam, Tajwid is a correct and regular method or way of reading the Al-Quran according to its makhraj, thickness and density of sound, long and short duration, whether it is resonant or not, as well as rhythm, tone and use of punctuation [1].

However, the reality in Indonesia shows that the level of Quranic literacy is still a serious challenge. Based on a research report presented by the Jakarta Institute of Quranic Sciences (IIQ) at the DPR/MPR RI Building, it was found that of 3,111 respondents spread across various regions of Indonesia, as many as 72.25% were categorized as having sufficient and insufficient Quranic reading skills [2]. The assessment was based on four main parameters in the science of tajwid, which reflects the great need for learning media that can bridge the Quranic literacy gap more effectively and inclusively.

In addition, a study conducted by Helmalia showed that 24.8% of students in an elementary school had difficulty recognizing the *hijaiyah* letters [3]. Meanwhile, Nurhayati reported that 42.5% of early childhood children did not recognize the *hijaiyah* letters well. These findings further emphasize the urgency of developing voice-based learning technology that is able to help people recognize, understand, and pronounce the *hijaiyah* letters interactively and adaptively, especially among children and beginners [4]. These findings further emphasize the urgency of developing voice-based learning technology that is able to help people recognize, understand, and pronounce the *hijaiyah* letters interactively and adaptively, especially among children and beginners.

In recent decades, advances in artificial intelligence (AI), particularly in speech signal processing, have enabled the creation of systems capable of recognizing and analyzing human voices with a high degree of accuracy. One approach that has proven effective in speech recognition tasks is the use of the **Convolutional Recurrent Neural Network (CRNN) algorithm**, which is a combination of a Convolutional Neural Network (CNN) for spatial feature extraction and a Recurrent Neural Network (RNN), specifically LSTM, for recognizing temporal patterns in speech data [5], [6], [7], [8]. On the other hand, the audio feature extraction process is also an important component that plays a role in improving the accuracy of voice classification. **Mel-Frequency Cepstral Coefficients (MFCC)** is one of the most widely used feature extraction methods in speech recognition systems due to its ability to effectively represent the characteristics of human voice [9], [10].

Research and development of applications utilizing the CRNN algorithm and MFCC feature extraction in speech recognition have been conducted in various fields, such as music genre classification, bird species identification, and voice-based attendance systems. However, its application in the context of Quranic learning, specifically for classifying the pronunciation of the *Hijaiyah* letters, has been limited. is still relatively rare. Therefore, this article was compiled as part of a literature review aimed at reviewing various relevant previous studies, understanding the methods used, and evaluating the effectiveness of technological approaches in supporting sound-based learning of the *hijaiyah* alphabet.

Thus, this article is expected to provide a strong theoretical and practical foundation for the development of educational applications that are able to increase Al-Quran reading literacy more widely, efficiently, and easily accessible to the general public.
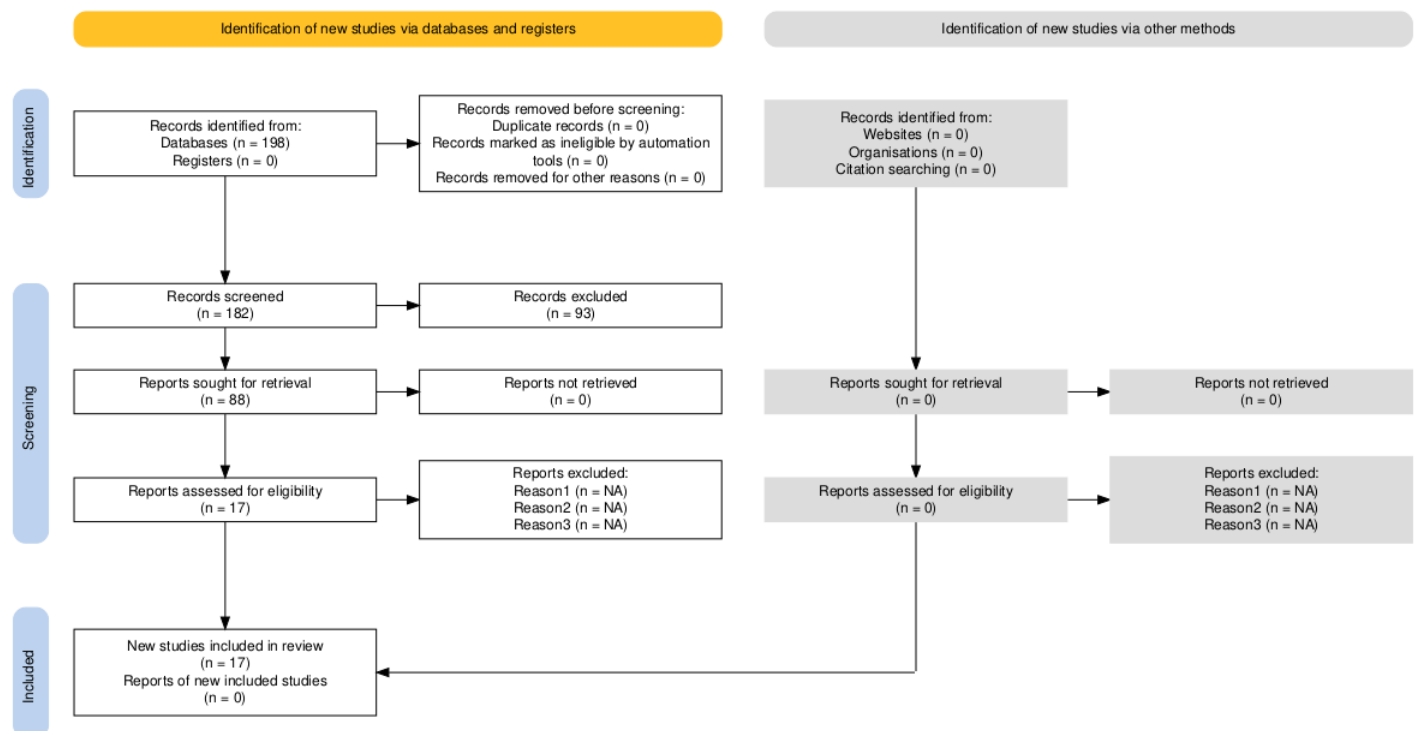


Fig. 1. PRISMA Diagram of the Study Selection Process

## II. RESEARCH METHODS

This study used a **Systematic Literature Review (SLR) approach** guided by **the PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) method** [11]. This method was used to ensure that the literature identification, selection, and evaluation process was carried out systematically, transparently, and could be replicated.
The steps in the process include:
1. **Literature Identification**

Literature was collected from relevant electronic database sources, including journals and conference proceedings focused on speech recognition, audio classification, and technology-based learning. In this phase, **198 documents** were identified from academic databases, with no documents from registers, websites, organizations, or other citation searches.
2. **Screening**
Of the initial 198 documents, screening was performed to remove duplicates and irrelevant studies. No duplicates or documents required removal due to automation were found, so all proceeded to the manual

screening stage. A total of **182 documents** were screened based on title and abstract, and **93** were eliminated for not matching the research focus, leaving **88 documents** for further review.

3. **Eligibility**

Of the 88 documents, only **17 reports** were fully accessible and met the inclusion criteria. Selection criteria included the use of CRNN, MFCC, or similar speech recognition algorithms.

4. **Inclusion**

All 17 reports that passed the eligibility stage were then included in the main review. No additional reports were included through other methods such as organization or citation searches.

By following this PRISMA flow, the literature study conducted is expected to cover previous research objectively, measurably, and relevantly to the focus of the study, namely the classification of the pronunciation of the *hijaiyah* letters. Voice-based.

## III.       RESULT AND DISCUSSION

### A.  Analysis of MFCC Feature Extraction in Arabic Voice Applications

Mel Frequency Cepstral Coefficient (MFCC) extraction is a very popular feature extraction method in speech signal processing. MFCC mimics human auditory perception of the sound frequency spectrum, making it suitable for languages with unique phonemes, such as Arabic.

Several studies have demonstrated the effectiveness of MFCC in various contexts. Putra applied MFCC and Naive Bayes to classify *cucak* bird sounds [12], while Hisyam and Prasetio [3] demonstrated that MFCC can be used effectively in vehicle siren classification with an accuracy of up to 94%. In the context of recognizing the laws of recitation of the Qur'an, which are closely related to Arabic phonology, Putra et al. utilized MFCC and HMM with an average accuracy of 80%, indicating that the features extracted by MFCC are powerful enough to distinguish phonetic patterns in Arabic [13].

From a cultural education perspective, Dewi and Suhartana [14] used MFCC to identify Balinese drum tones. This shows that MFCC is quite flexible when applied to distinctive sounds, including the possible sounds of the *hijaiyah* letters.

Research by Nursholihatun et al. also shows that the use of MFCC for speaker identification can achieve an accuracy of up to 100% on training data and remains high even with the addition of noise (SNR), demonstrating the power of MFCC in diverse environmental conditions [14]. This is an important indication for Arabic speech recognition applications in learning or classifying *hijaiyah* letter sounds which are usually pronounced in various acoustic conditions.

### B.  Application of CRNN in Processing of Tempo-Sound Data

A Convolutional Recurrent Neural Network (CRNN) is a combination of CNN and RNN, capable of simultaneously capturing spatial and temporal features in audio data. CNNs are used to extract features from spectral representations (e.g., spectrograms), while RNNs, specifically LSTMs or GRUs, capture the temporal dynamics of the data sequence.

Yuslan and Adiwijaya proved the effectiveness of CRNN in hadith text classification with an accuracy of 80.79% [15], outperforming CNN or RNN used separately. Fatichin introduced a variant of dilated-CRNN with MFCC and data augmentation in music genre classification, which significantly outperformed standard CRNN [16].

With its layered structure that can handle long-term temporal relationships, CRNN is ideal for analyzing tempo-driven sound data such as Arabic letter pronunciations that have specific length-shortness dynamics and stresses (e.g., mad, sukun, or tasydid). Therefore, CRNN can be a relevant solution for *hijaiyah* sound classification because it is able to capture complex transitions between phonemes.

### C.  Potential Combination of CRNN and MFCC for Classification of Hijaiyah Letters

Based on previous analysis, the combination of MFCC and CRNN shows great potential in developing a speech-based *hijaiyah* letter classification system. MFCC acts as a feature extraction stage that filters important characteristics from the voice, while CRNN serves as a robust classification model for processing the temporal sequence of the signal.

This combination has been proven successful in other domains. Fatichin showed that dilated CRNN and MFCC outperformed other methods in music genre classification [16]. Similarly, Ihsanti et al. used CNN with MFCC for owl call classification and achieved near-perfect accuracy of 99.8% [17]. A similar approach can be applied to the sounds of the *hijaiyah* letters which have distinctive phonetic and rhythmic patterns.

Considering the complexity of the pronunciation of the *hijaiyah* letters, such as the hamzah, ha ( ح ), or 'ain ( ع ) sounds, a classification system requires a fine-grained feature representation as well as a model capable of capturing temporal dynamics. Therefore, the application of MFCC as an input feature and CRNN as a model architecture is very suitable for recognizing or classifying *hijaiyah* letter sounds in educational, religious, or rehabilitative contexts.

Table 1. Review Result

| No | Researchers | Method/Technology | Main Results & Applications |
|---|---|---|---|
| 1. | Muhammad Yuslan Abu Bakar, Adiwijaya [15] | CRNN (combination of CNN & RNN) | 80.79% accuracy in hadith text classification; CRNN is better than single CNN/RNN |
| 2. | Soraya Ihsanti, Wikky F. Al Maki [17] | CNN + MFCC | 99.8% accuracy in owl sound classification |

| No | Researchers | Method/Technology | Main Results & Applications |
|---|---|---|---|
| 3. | Mochammad R. Fatichin. [16] | Dilated-CRNN + MFCC + data augmentation | Outperforming standard CRNN for music genre classification |
| 4. | Muhammad R. Putra [12] | MFCC + Naive Bayes | Accuracy of 52–90% in classifying similar cucak bird sounds |
| 5. | M. Ibn Hisham, BH Prasetio [3] | MFCC + CNN | 94% accuracy in classifying priority vehicle siren sounds |
| 6. | Faisal D. Adhinata [18] | MFCC + GMM | 81.18% accuracy in gender classification based on voice |
| 7. | Muhammad Afif Ma'ruf [19] | MFCC + LVQ | Accuracy of 90% (without mask) and 86.2% (with mask) for voice attendance verification |
| 8. | Fuad F. Surenggana [20] | MFCC + KNN | 85.5% accuracy in music mood classification (happy, sad, angry, relax) |
| 9. | Oddy V. Putra [13] | MFCC + HMM | Average accuracy of 80% in the classification of the laws of mad reading of the Qur'an |
| 10. | Ni Kadek Yulia Dewi, I Ketut Gede Suhartana [21] | MFCC + KNN | 90% accuracy in identifying the basic tone of Balinese drums, used for application-based traditional music training |
| 11. | Galih Ajinurseto, La Ode Bakrim, Nur Islamuddin [22] | MFCC (no specific classification) | Desktop-based speech recognition system, 90% accuracy under ideal conditions and 76.6% under non-ideal conditions |
| 12. | Tobias Sion Julian, Fitri Utaminingrum, Dahnial Syauqy [23] | MFCC + CNN | Voice command system for smart wheelchairs; 85% accuracy for 4 voice commands ("Go", "Left", "Right", "Stop") based on Jetson TX2 |
| 13. | Ricky Aurelius N. Diaz, Ni Luh Gede Pivin, Komang Budiarta [24] | MFCC vs Spectral + KNN | MFCC outperforms spectral in gender-based voice classification; training accuracy is 84.18% and testing accuracy is 74.71%. |
| 14. | Muhammad Nabil Aljufri, Barlian H. Prasetio [25] | MFCC + Artificial Neural Network | Raspberry Pi-based voice stress detection; 90% accuracy in testing, and 76% in the final ANN model |
| 15. | Andre Julio S. Marbun, Heriyanto, Frans R. Kodong [26] | MFCC + SVM | Accuracy of 85.71% (male) and 92.21% (female) in voice emotion classification using the RAVDESS dataset |
| 16. | Femmy Fauziah, Iwan I. Tritoasmoro, Syamsul Rizal [27] | MFCC + Vector Quantization (VQ) | Voice recognition system for door security; 88% accuracy (quiet, 10 cm) and 68% (noisy, 15 cm) using Arduino and Python |
| 17. | Erina Nursholihatun, Sudi MA Sasongko, Abdullah Zainuddin [14] | MFCC + Backpropagation Artificial Neural Network | Speech recognition accuracy reaches 100% on training data, 86% on noise-free test data, and increases from 45% to 92% depending on the SNR level. 100% rejection of out-of-database speakers. Used for a voice-based speaker identification system for 22–24-year-old males using MATLAB. |

Based on the results of a literature study of 17 relevant publications, it can be concluded that the combination of the Mel-Frequency Cepstral Coefficients (MFCC) and Convolutional Recurrent Neural Network (CRNN) algorithms shows high potential in sound classification tasks, especially for the pronunciation of *hijaiyah* letters.

1. MFCC has been proven to be effective in efficiently extracting acoustic features that represent the phonetic characteristics of human voices. This has been demonstrated in various domains, ranging from the classification of bird calls, sirens, emotions, to the recognition of the rules of recitation of the Qur'an. MFCC has robustness to noise and flexibility in different acoustic conditions, making it very suitable for Arabic which is rich in phonemes and intonation [3], [13], [14].

2. CRNN, as a combination of CNN and RNN, is able to capture spatial and temporal patterns in sound data simultaneously. This advantage makes it superior in managing tempo-driven audio signals such as the pronunciation of *hijaiyah* letters, which contain elements of length and shortness, stress, and transitions between phonemes [15], [16].

3. The combination of MFCC and CRNN has successfully increased classification accuracy in various studies—including on music genres and hadith texts—and is considered to have higher accuracy than conventional methods such as CNN, RNN, or other classification algorithms [15], [16], [17].

4. This study strengthens the argument that the development of a speech-based Quran learning system, specifically to assist in the classification and pronunciation training of the *Hijaiyah* letters, can be significantly enhanced by the application of MFCC + CRNN technology. This presents a significant opportunity to support the improvement of Quran literacy in Indonesia, particularly among children, beginners, and the general public.

Thus, the results of this study provide a strong theoretical basis for the development of artificial intelligence-based educational applications for automatic, adaptive, and accurate classification of *hijaiyah* letter pronunciation.

## V.     CONCLUSION

This study concludes that the combination of Mel-Frequency Cepstral Coefficients (MFCC) and Convolutional Recurrent Neural Network (CRNN) offers a highly effective framework for classifying the pronunciation of *hijaiyah* letters, providing a promising foundation for AI-driven Qur'anic literacy tools. MFCC proves to be a robust feature extraction technique suitable for the rich phonetic characteristics of Arabic sounds, while CRNN excels at capturing the temporal dynamics essential to accurate pronunciation modeling. Together, they form a synergistic system well-suited for developing adaptive and interactive tajwid learning applications. As a direction for future work, researchers are encouraged to implement and evaluate this model using real-world *hijaiyah* pronunciation datasets

across varied age groups and accents. Additionally, integrating user feedback mechanisms and gamification could enhance learner engagement and make AI-based Qur'anic education more accessible, especially in underserved or remote communities.

## REFERENCES

[1] I. F. Muslim, S. Ranam, and P. Priyono, "Peningkatan Kemampuan Membaca Alquran dengan Pelatihan," *PUNDIMAS: Publikasi Kegiatan Abdimas*, vol. 1, no. 2, pp. 70–73, Jun. 2022, doi: 10.37010/pnd.v1i2.680.

[2] Institut Ilmu Al-Qur'an, "Tim IIQ Jakarta Paparkan Hasil Riset Tingginya Buta Aksara Al-Qur'an di Gedung DPR-MPR RI Senayan," Institut Ilmu Al-Qur'an Website. Accessed: Jun. 02, 2025. [Online]. Available: https://iiq.ac.id/berita/tim-iiq-jakarta-paparkan-hasil-riset-tingginya-buta-aksara-al-quran-di-gedung-dpr-mpr-ri-senayan/

[3] M. I. Hisyam and B. H. Prasetio, "Klasifikasi Suara Sirene Kendaraan berbasis MFCC untuk Meningkatkan Efisiensi Sistem Keamanan Lalu Lintas," *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, vol. 8, no. 7, 2024.

[4] S. Susanti and S. Nurhayati, "Penerapan Metode Iqro'dalam Mengenalkan Huruf Hijaiyah Pada Anak Usia Dini," *WALADUNA: Jurnal Pendidikan Islam Anak Usia Dini*, vol. 5, no. 2, pp. 13–23, 2022, doi: https://doi.org/10.12928/waladuna.v5i2.533.

[5] U. Ayvaz, H. Gürüler, F. Khan, N. Ahmed, T. Whangbo, and A. Akmalbek Bobomirzaevich, "Automatic Speaker Recognition Using Mel-Frequency Cepstral Coefficients Through Machine Learning," *Computers, Materials & Continua*, vol. 71, no. 3, pp. 5511–5521, 2022, doi: 10.32604/cmc.2022.023278.

[6] Z. Kh. Abdul and A. K. Al-Talabani, "Mel Frequency Cepstral Coefficient and its Applications: A Review," *IEEE Access*, vol. 10, pp. 122136–122158, 2022, doi: 10.1109/ACCESS.2022.3223444.

[7] N. R. Haddaway, M. J. Page, C. C. Pritchard, and L. A. McGuinness, "PRISMA2020: An R package and Shiny app for producing PRISMA 2020-compliant flow diagrams, with interactivity for optimised digital transparency and Open Synthesis," *Campbell Systematic Reviews*, vol. 18, no. 2, Jun. 2022, doi: 10.1002/cl2.1230.

[8] M. R. Putra, F. Nurdiyansyah, and A. Y. Rahman, "Klasifikasi Jenis Burung Cucak Berdasarkan Suara Menggunakan MFCC Dan Naive Bayes," *JURNAL FASILKOM*, vol. 14, no. 2, pp. 463–470, 2024.

[9] O. V. Putra, Faisal Reza Pradana, and Jordan Istiqlal Qalbi Adiba, "Mad Reading Law Classification Using Mel Frequency Cepstal Coefficient (MFCC) and Hidden Markov Model (HMM)," *Procedia of Engineering and Life Science*, vol. 2, Dec. 2021, doi: 10.21070/pels.v2i0.1148.

[10] E. Nursholihatun, S. M. Sasongko, and A. Zainuddin, "Identifikasi Suara Menggunakan Metode Mel Frequency Cepstrum Coefficients (MFCC) dan Jaringan Syaraf Tiruan Backpropagation," *DIELEKTRIKA*, vol. 7, no. 1, p. 48, Feb. 2020, doi: 10.29303/dielektrika.v7i1.232.

[11] M. Y. Abu Bakar and A. Adiwijaya, "Klasifikasi Teks Hadis Bukhari Terjemahan Indonesia Menggunakan Recurrent Convolutional Neural Network (CRNN)," *Jurnal Teknologi Informasi dan Ilmu Komputer*, vol. 8, no. 5, p. 907, Oct. 2021, doi: 10.25126/jtiik.2021853750.

[12] M. R. Fatichin, A. R. Hermawan, R. A. S. Siahaan, and R. Indraswari, "Dilated-Convolutional Recurent Neural Network untuk Klasifikasi Genre Musik," *Jurnal Teknik Informatika dan Sistem Informasi*, vol. 10, no. 3, pp. 439–448, 2024, doi: https://doi.org/10.28932/jutisi.v10i3.9347.

[13] S. Ihsanti and W. F. Al Maki, "Klasifikasi Genus Burung Hantu Berdasarkan Suara Menggunakan Convolutional Neural Network," *eProceedings of Engineering*, vol. 11, no. 4, 2024.

[14] F. D. Adhinata, D. P. Rakhmadani, and A. J. T. Segara, "Pengenalan Jenis Kelamin Manusia Berbasis Suara Menggunakan MFCC dan GMM," *Journal of Dinda: Data Science, Information Technology, and Data Analytics*, vol. 1, no. 1, pp. 28–33, 2021.

[15] M. A. Ma'ruf, A. Aranta, and F. Bimantoro, "Verifikasi Suara Mahasiswa Sebagai Alternatif Presensi Kehadiran Menggunakan Ekstraksi Fitur MFCC Dan Klasifikasi LVQ," *Jurnal Teknologi Informasi, Komputer, dan Aplikasinya (JTIKA)*, vol. 4, no. 2, pp. 171–181, 2022, doi: https://doi.org/10.29303/jtika.v4i2.211.

[16] F. F. Surenggana, A. Aranta, and F. Bimantoro, "Klasifikasi Mood Musik Menggunakan K-Nearest Neighbor dan Mel Frequency Cepstral Coefficients," *Jurnal Teknologi Informasi, Komputer, dan Aplikasinya (JTIKA)*, vol. 4, no. 2, pp. 263–276, 2022, doi: https://doi.org/10.29303/jtika.v4i2.191.

[17] N. K. Y. Dewia and I. K. G. Suhartanaa, "Idenfikasi Nada Dasar Kendang Menggunakan MFCC dan KNN," *J. Elektron. Ilmu Komput. Udayana,* vol. 11, no. 4, pp. 797–802, 2023.

[18] G. Ajinurseto and N. Islamuddin, "Penerapan metode mel frequency cepstral coefficients pada sistem pengenalan suara berbasis desktop," *Infomatek*, vol. 25, no. 1, pp. 11–20, 2023.

[19] T. S. Julian, F. Utaminingrum, and D. Syauqy, "Sistem Voice Command pada Kursi Roda Pintar menggunakan MFCC dan CNN berbasis Jetson TX2," *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, vol. 6, no. 11, pp. 5505–5510, 2022.

[20] R. A. N. Diaz, N. L. G. P. Suwirmayanti, and K. Budiarta, "Perbandingan Kualitas Pengenalan Suara untuk Ekstraksi Fitur Menggunakan MFCC dan Spectral," *Naratif : Jurnal Nasional Riset, Aplikasi dan Teknik Informatika*, vol. 6, no. 1, pp. 58–63, Jul. 2024, doi: 10.53580/naratif.v6i1.281.

[21] M. N. Aljufri and B. H. Prasetio, "Sistem Deteksi Tingkat Stress Menggunakan Suara dengan Metode Jaringan Saraf Tiruan dan Ekstraksi Fitur MFCC berbasis Raspberry Pi," *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, vol. 6, no. 11, pp. 5278–5285, 2022.

[22] A. J. Marbun and F. R. Kodong, "Implementation of Mel-Frequency Cepstral Coefficient As Feature Extraction Method On Speech Audio Data," *Telematika: Jurnal Informatika dan Teknologi Informasi*, vol. 8, no. 36, pp. 11839–11848, 2024, doi: https://doi.org/10.31315/telematika.v21i3.12339.

[23] F. Fauziah, I. I. Tritoasmoro, and S. Rizal, "Sistem Keamanan Berbasis Pengenalan Suara Sebagai Pengakses Pintu Menggunakan Metode Mel Frequency Cepstral Coefficient (MFCC)," *eProceedings of Engineering*, vol. 8, no. 6, 2021.