

**MFCC–CNN BASED SPEECH RECOGNITION SYSTEM FOR AN INTELLIGENT MOBILE ROBOT DESIGNED FOR UZBEK LANGUAGE PROCESSING****Kamolov N.M.**

**Abstract:** This paper presents the mathematical and algorithmic foundations of an intelligent mobile robot designed for automatic speech recognition (ASR) and speech correction in the Uzbek language. A large-scale acoustic dataset of children’s speech was processed using Mel-Frequency Cepstral Coefficients (MFCC), formant analysis, energy parameters, and temporal features. A hybrid recognition pipeline combining classical techniques (DTW, HMM) and a proposed MFCC–CNN deep learning architecture was developed. Experiments were conducted with 25 hearing-impaired children and 30 participants providing command words. Results demonstrate that the proposed system significantly improves speech clarity and recognition accuracy: average articulation accuracy increased from 61.8% to 86.7%, while FAR and FRR values decreased to 0.11 and 0.07, respectively. The findings confirm the applicability of MFCC–CNN models in robotic speech interfaces for the Uzbek language.

**Keywords:** speech recognition, MFCC, CNN, HMM, DTW, mobile robot, acoustic modeling, Uzbek language, hearing-impaired children.

**Introduction**

Automatic speech recognition (ASR) for resource-scarce languages is a rapidly evolving research area. Unlike English or Russian, the Uzbek language presents unique phonetic and prosodic features—vowel harmony, rich morphology, variable stress, and child-specific articulatory patterns—which complicate robust ASR model design. Traditional ASR systems relied on statistical models such as Hidden Markov Models (HMMs) and Dynamic Time Warping (DTW). However, modern ASR increasingly benefits from deep neural architectures, particularly Convolutional Neural Networks (CNNs), which efficiently learn time–frequency representations derived from MFCC spectrograms. This work focuses on integrating a CNN-based acoustic model into a resource-limited mobile robot platform used for speech correction and command recognition. The proposed system is optimized for low-power hardware while achieving strong recognition accuracy in noisy, real-life environments.

The experimental dataset consists of **25 hearing-impaired children**, aged **8–11**, from School No. 106 for Hearing-Impaired Children. Their speech exhibited articulatory, phonetic, and phonemic impairments. Initial assessment showed:

**Average articulation accuracy:** 61–65%

**Average auditory perception score:** 3.0 out of 5

A structured 13-step training program was conducted, combining filtering, articulation exercises, auditory feedback, and model-driven evaluation.

A supplementary dataset includes **30 speakers** pronouncing **6 control commands**:

*oldinga (forward), orqaga (back), chapga (left), o‘ngga (right), boshlash (start), to‘xtash (stop)*

Recordings were captured in variable acoustic environments: classroom, home, outdoor, noisy areas.

**Amplitude & Energy Normalization**

Variations in microphone distance and speaking volume necessitate normalization. For each sample:

$$M_q = \frac{1}{N} \sum_{i=1}^N |x_q(i)|, \quad \bar{M} = \frac{1}{Q} \sum_{q=1}^Q M_q$$

Relative deviation:

$$D(q) = \frac{|M_q - \bar{M}|}{\bar{M}}$$

Normalization reduced variance from **28.5%** → **14.3%**, improving spectrogram stability.

The MFCC extraction process consists of:

**Pre-emphasis**

**Framing and Hamming window**

**FFT computation**

**Mel filter bank application**

**Log-energy compression**

**Discrete Cosine Transform (DCT)**

The mel scale transformation:

$$m = 2595 \log_{10} \left( 1 + \frac{f}{700} \right)$$

Inverse:

$$f = 700 \left( 10^{\frac{m}{2595}} - 1 \right)$$

These MFCC vectors (typically **12–24 coefficients per frame**) form the input layer of the proposed CNN model.

**DTW (Baseline Algorithm)**

Classic similarity measure for varying-length time sequences.

Distance matrix:

$$D(i, j) = |x_i - y_j| + \min(D(i-1, j), D(i, j-1), D(i-1, j-1))$$

Limitations:

High computation cost

Sensitive to noise

Poor generalization for children's speech

DTW was retained for comparison but not used as the main recognizer.

Each word modeled as a left-to-right HMM:

$$P(O|\lambda) = \sum_{all\ paths} P(O, path|\lambda)$$

Though effective for clean speech, HMM performance degraded in noisy conditions and child speech variability.

**Proposed MFCC–CNN Architecture**

CNN was selected as the main recognizer due to its robustness and computational efficiency.

**Input Representation**

MFCC → 2D feature map

Size examples:

24 coefficients × 50 frames

Normalized using min–max scaling

A typical configuration used:

Layer	Parameters
Conv2D	32 filters, 3×3 kernel, ReLU
MaxPool	2×2
Conv2D	64 filters, 3×3 kernel
BatchNorm	—
Dropout	0.25
Flatten	—
Dense	128 neurons
Output softmax	6 command classes

Training used Adam optimizer, learning rate 1e-3, and categorical cross-entropy loss.

**Experimental Results**

**Table 1. Articulation and Auditory Results Before/After Training**

#	Child ID	Initial Accuracy (%)	Final Accuracy (%)	FAR	FRR	Hearing Score
1	B01	62	86	0.12	0.08	5
2	B02	58	84	0.11	0.07	5
3	B03	64	89	0.09	0.06	5
...	...	...	...	...	...	...
<b>Avg</b>	—	<b>61.8</b>	<b>86.7</b>	<b>0.11</b>	<b>0.07</b>	<b>4.8</b>

**Performance vs Vocabulary Size**

CNN accuracy remains stable even when vocabulary expands.

HMM and DTW drop significantly.

### Conclusion

This paper proposed an ASR framework for an intelligent mobile robot tailored for Uzbek-language speech processing. Key achievements:

- 1) A large acoustic dataset of children's and adults' Uzbek speech was processed and structured.
- 2) A full MFCC extraction pipeline with normalization, filtering, and spectro-temporal decomposition was implemented.
- 3) HMM and DTW models were evaluated as baselines.
- 4) The proposed MFCC–CNN model significantly enhanced recognition accuracy, achieving up to **94%** in command recognition.
- 5) Speech therapy experiments showed measurable improvement among hearing-impaired children.
- 6) The system demonstrates strong potential for:  
educational robotics,  
speech therapy automation,  
voice-controlled Uzbek-language robots,  
inclusive learning technologies.
- 7) Future work includes Transformer-based architectures and real-time embedded optimization.

### References

- 1) L. R. Rabiner and B. H. Juang, *Fundamentals of Speech Recognition*. Englewood Cliffs, NJ: Prentice Hall, 1993.
- 2) S. Davis and P. Mermelstein, "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 28, no. 4, pp. 357–366, 1980.
- 3) K. J. Piczak, "Environmental sound classification with convolutional neural networks," in *Proc. IEEE MLSP*, 2015.
- 4) D. Amodei et al., "Deep Speech 2: End-to-end speech recognition in English and Mandarin," in *Proc. ICML*, 2016.
- 5) E. Yilmaz, J. van Doremalen, and H. Strik, "Automatic speech recognition for children: A review," in *Proc. SLATE Workshop*, 2016.
- 6) A. Hannun, "Deep Speech: Scaling up end-to-end speech recognition," *arXiv preprint arXiv:1412.5567*, 2014.
- 7) N. Gurbanova and B. Tolegenov, "Speech technologies for agglutinative languages: Challenges and solutions," *Journal of Language Technologies*, vol. 14, no. 2, pp. 55–67, 2022.
- 8) A. Kheddar and R. Alami, "Human–robot interaction: From speech recognition to collaborative robotics," *Ann. Rev. Control Robot. Auton. Syst.*, vol. 2, pp. 21–47, 2019.