

Speech Emotion Recognition using Fuzzy-Based CNN with Optimal Feature Selection using Beluga Whale Optimization

Chandupatla Deepika
KLEF, KL University AP

Dr. K. Swarna

Associate Professor, Department of CSE, KLEF.

Dr. Syed Khasim

Dr.Samuel George Institute of Engineering & Technology, Markapur, Prakasam DT. Andhra Pradesh, India.

Abstract

Speech Emotion Recognition (SER) plays a vital role in human-computer interaction, healthcare, security, and psychological analysis. However, existing SER systems face challenges due to noise, speaker variations, and the overlapping nature of emotional speech features. To address these complexities, this thesis proposes a Fuzzy-Based Convolutional Neural Network (Fuzzy-CNN) integrated with Beluga Whale Optimization (BWO) for optimal feature selection. The proposed system consists of four major stages: speech pre-processing, feature extraction, fuzzy-CNN classification, and BWO-driven feature selection. Audio features such as MFCC, Delta MFCC, spectral features, and Chroma vectors are extracted from speech signals. A fuzzy layer is embedded within the CNN to manage ambiguity in emotional speech patterns, enabling better decision-making under uncertainty. To reduce redundancy and enhance model accuracy, the Beluga Whale Optimization algorithm is utilized to select the most relevant features. Experimental evaluations on benchmark datasets demonstrate that the proposed hybrid approach significantly improves recognition accuracy, reduces computational complexity, and performs better than traditional machine learning and deep learning methods. Comparative analysis shows enhanced precision, recall, F1-score, and robustness under noisy conditions. This research contributes an interpretable fuzzy-deep learning framework and introduces BWO-based optimization to advance the effectiveness of SER systems. The findings highlight the potential of hybrid intelligent models for real-time emotion-aware applications.

Keywords: Speech emotion recognition, Fuzzy neural networks, Convolutional neural networks, Beluga whale optimization, Feature selection, MFCC, Deep learning

1. INTRODUCTION

Human emotions are expressed through multiple channels including facial expressions, body language, and vocal cues. Among these, speech carries rich emotional information that reflects a person's psychological state, intentions, and reactions. The ability to automatically recognize emotions from speech has become increasingly important as technology integrates more deeply into daily life. Applications span diverse domains from mental health monitoring systems that detect depression or anxiety through voice analysis, to customer service platforms that adapt responses based on caller emotions, to smart home assistants that respond empathetically to user frustration or excitement (Akçay and Oğuz, 2020).

10.48047/jocaaa.2024.33.05.40

Despite significant progress in speech processing and artificial intelligence, accurate emotion recognition from speech remains challenging. Unlike visual emotion recognition where facial expressions provide relatively clear indicators, speech-based emotions are subtle and complex. The same words spoken with different intonations, pitch variations, and speaking rates can convey completely different emotional states. Moreover, individual differences in speaking styles, cultural backgrounds, and linguistic patterns add layers of complexity that make universal emotion recognition difficult (Lalitha et al., 2015).

Current speech emotion recognition systems primarily rely on traditional machine learning classifiers like Support Vector Machines, Random Forests, and Hidden Markov Models, or more recently, deep learning approaches including standard CNNs and Recurrent Neural Networks. However, these methods face several limitations. Traditional classifiers struggle with high-dimensional feature spaces and require extensive manual feature engineering. Standard deep learning models, while powerful, often treat emotional categories as crisp boundaries and fail to capture the fuzzy, ambiguous nature of human emotions. For instance, distinguishing between anger and frustration or between happiness and excitement involves subtle differences that exist on a continuum rather than as discrete categories (Kwon, 2020).

The feature extraction and selection process presents another major challenge. Speech signals can yield hundreds of potential features including temporal, spectral, prosodic, and voice quality characteristics. Many of these features are redundant or irrelevant for emotion recognition, yet selecting the optimal subset manually is impractical. Including all features increases computational complexity and risks overfitting, while naive feature reduction may discard important emotional indicators. This challenge necessitates intelligent optimization approaches that can automatically identify the most discriminative features for emotion classification (Zhang et al., 2021).

This research addresses these gaps by proposing a novel hybrid architecture that combines three powerful concepts: fuzzy logic for handling emotional ambiguity, convolutional neural networks for automatic feature learning, and Beluga Whale Optimization for intelligent feature selection. The fuzzy component acknowledges that emotions exist on continuums with overlapping boundaries. The CNN architecture automatically learns hierarchical representations from speech features. The BWO algorithm, inspired by the social foraging behavior of beluga whales, efficiently searches the feature space to identify optimal subsets that maximize classification performance while minimizing redundancy.

The research is guided by three primary questions. First, how can fuzzy logic be effectively integrated into CNN architectures to better handle the inherent ambiguity in emotional speech? Second, can nature-inspired optimization algorithms like BWO outperform traditional feature selection methods in identifying relevant speech features for emotion recognition? Third, does the proposed hybrid system demonstrate measurable improvements in accuracy, computational efficiency, and robustness compared to existing approaches?

This paper makes several key contributions to the field of speech emotion recognition. We present a novel fuzzy-CNN architecture specifically designed for emotion classification that outperforms conventional approaches. We introduce the application of Beluga Whale Optimization to speech feature selection, demonstrating its effectiveness in this domain. We provide comprehensive experimental validation on benchmark datasets showing improvements in multiple performance metrics. Finally, we

offer practical insights for implementing hybrid intelligent systems in real-world emotion-aware applications.

The remainder of this paper is organized as follows. Section 2 outlines the research objectives and scope. Section 3 reviews relevant literature on speech emotion recognition, fuzzy systems, CNNs, and optimization algorithms. Section 4 details the proposed methodology including system architecture and implementation. Section 5 presents experimental results and comparative analysis. Section 6 discusses findings, implications, and limitations. Section 7 concludes with contributions and future research directions.

2. OBJECTIVES

The specific objectives guiding this research are:

- **Primary Objective:** To develop and validate a hybrid Fuzzy-CNN architecture integrated with Beluga Whale Optimization that achieves superior emotion recognition accuracy compared to baseline machine learning and deep learning methods, with measurable improvements of at least 8-12% in classification metrics on standard benchmark datasets.
- **Secondary Objective 1:** To design an effective fuzzy inference layer that can be seamlessly integrated within CNN architectures to handle uncertainty and overlapping boundaries in emotional speech features, quantifying the contribution of fuzzification to overall system performance.
- **Secondary Objective 2:** To implement and optimize the Beluga Whale Optimization algorithm for feature selection in speech emotion recognition, demonstrating its ability to reduce feature dimensionality by 40-60% while maintaining or improving classification accuracy compared to using complete feature sets.
- **Secondary Objective 3:** To conduct comprehensive experimental evaluation across multiple emotion datasets and noise conditions, establishing the robustness and generalizability of the proposed approach through statistical validation and comparative analysis against at least four baseline methods.
- **Secondary Objective 4:** To analyze computational complexity and real-time performance characteristics of the proposed system, providing practical guidelines for deployment in resource-constrained environments and real-time emotion-aware applications.

3. SCOPE OF STUDY

This research operates within the following defined boundaries:

Temporal and Data Scope:

- Focus on publicly available benchmark speech emotion datasets recorded between 2015-2024
- Emphasis on acted emotional speech databases with clear ground truth labels
- Consideration of both speaker-dependent and speaker-independent recognition scenarios
- Exclusion of spontaneous or naturalistic emotional speech which introduces additional labeling ambiguities

Emotional Categories:

- Primary focus on six basic emotions: anger, disgust, fear, happiness, sadness, and neutral
- Exclusion of complex or secondary emotions like pride, shame, or anticipation
- Recognition at utterance level rather than continuous emotion tracking within speech

Technical Boundaries:

- Implementation using Python-based deep learning frameworks suitable for reproducibility
- Feature extraction limited to acoustic and prosodic features, excluding linguistic content analysis
- CNN architectures constrained to 2D convolutions operating on spectrogram-like representations
- Optimization limited to feature selection rather than hyperparameter tuning or architecture search

Language and Cultural Scope:

- Primary validation on English language datasets
- Acknowledgment that cross-cultural emotion expression varies but detailed cultural analysis is beyond scope
- Framework designed to be language-agnostic though validation focuses on English corpora

Application Context:

- System designed for offline batch processing and near-real-time applications •
- Computational requirements suitable for modern GPU-equipped servers
- Exclusion of extremely resource-constrained edge devices like microcontrollers
- Privacy considerations addressed through local processing capability but detailed security analysis excluded

4. LITERATURE REVIEW

4.1 Foundations of Speech Emotion Recognition

Speech emotion recognition has evolved significantly over the past two decades, progressing from simple acoustic analysis to sophisticated deep learning systems. Early SER research focused on identifying prosodic features like pitch, intensity, and speaking rate that correlate with different emotional states (Schuller et al., 2003). These handcrafted features were fed into traditional classifiers including Gaussian Mixture Models, Support Vector Machines, and Decision Trees.

The introduction of Mel-Frequency Cepstral Coefficients revolutionized speech processing by providing compact representations that capture important perceptual characteristics of audio signals (Muda et al., 2010). MFCCs and their temporal derivatives became standard features for SER systems. However, researchers quickly recognized that emotions manifest through multiple acoustic dimensions beyond MFCCs, leading to extensive feature extraction including spectral features, voice quality measures, and formant characteristics.

10.48047/jocaaa.2024.33.05.40

A persistent challenge in SER has been the subjective and context-dependent nature of emotions. The same utterance might be labeled differently by different annotators, and cultural factors influence both emotion expression and perception (Busso et al., 2008). This inherent ambiguity suggests that crisp classification boundaries may be inappropriate, motivating the exploration of fuzzy and probabilistic approaches.

4.2 Deep Learning for Speech Emotion Recognition

The deep learning revolution transformed speech emotion recognition by enabling automatic feature learning from raw or minimally processed audio. Convolutional Neural Networks, originally developed for image processing, proved effective when applied to spectrograms and other time-frequency representations of speech (Badshah et al., 2017). CNNs automatically learn hierarchical filters that capture local patterns and progressively build higher-level representations.

Recurrent Neural Networks, particularly Long Short-Term Memory networks, showed promise for modeling temporal dependencies in emotional speech (Mirsamadi et al., 2017). The sequential nature of speech makes RNNs a natural choice, allowing models to capture how emotions evolve across utterances. Hybrid CNN-LSTM architectures emerged that combine CNN's spatial feature learning with LSTM's temporal modeling capabilities.

Attention mechanisms further improved SER performance by allowing models to focus on emotionally salient portions of speech while ignoring irrelevant segments (Zhao et al., 2019). More recently, transformer architectures and self-supervised learning approaches have achieved state-of-the-art results by pre-training on large unlabeled speech corpora before fine-tuning for emotion classification.

Despite these advances, standard deep learning approaches have limitations. They typically require large labeled datasets which are expensive to create for emotional speech. Many architectures are black boxes that provide little interpretability about which acoustic characteristics drive predictions. Most critically for this research, conventional neural networks treat emotion categories as mutually exclusive classes with sharp boundaries, failing to capture emotional ambiguity.

4.3 Fuzzy Logic in Emotion Recognition

Fuzzy logic provides a mathematical framework for reasoning with uncertainty and imprecision, making it well-suited for emotion recognition where boundaries between categories are inherently vague (Zadeh, 1965). Several researchers have explored fuzzy approaches for SER, typically by applying fuzzy clustering or fuzzy rule-based systems to classify acoustic features (Wu et al., 2011).

Fuzzy membership functions can represent the degree to which speech characteristics belong to different emotional categories. For instance, a speech sample might have 0.7 membership in "anger" and 0.4 membership in "frustration," capturing the reality that these emotions share acoustic similarities. Fuzzy inference systems use linguistic rules like "if pitch is high and intensity is very high then emotion is likely anger" to make classification decisions.

10.48047/jocaaa.2024.33.05.40

However, traditional fuzzy systems still require manual design of membership functions and rules, limiting their scalability to complex, high-dimensional speech features. This motivated research into neuro-fuzzy systems that combine neural networks' learning capabilities with fuzzy logic's interpretability (Jang, 1993). Adaptive Neuro-Fuzzy Inference Systems automatically tune membership functions and rules through training data.

Recent work has begun exploring integration of fuzzy layers within deep neural networks, though application to speech emotion recognition remains limited (Liu and Wang, 2020). These fuzzy neural networks promise to combine automatic feature learning, uncertainty handling, and some degree of interpretability.

4.4 Feature Selection and Optimization

The curse of dimensionality poses significant challenges in speech emotion recognition where hundreds of features can be extracted from each utterance. Redundant and irrelevant features increase computational costs, risk overfitting, and can actually degrade classification performance (Guyon and Elisseeff, 2003). Feature selection aims to identify the minimal subset that retains maximal discriminative information.

Traditional feature selection methods include filter approaches that rank features based on statistical measures, wrapper methods that evaluate subsets using classifier performance, and embedded techniques that perform selection during model training (Saeys et al., 2007). Each approach has tradeoffs between computational efficiency and selection quality.

Nature-inspired metaheuristic algorithms have emerged as powerful tools for feature selection by treating it as an optimization problem. Genetic Algorithms, Particle Swarm Optimization, and Ant Colony Optimization have been successfully applied to SER feature selection (Kerkeni et al., 2019). These algorithms use mechanisms inspired by natural phenomena to efficiently search large feature spaces.

4.5 Beluga Whale Optimization Algorithm

Beluga Whale Optimization is a relatively recent metaheuristic algorithm inspired by the social behaviors of beluga whales during swimming, foraging, and predator avoidance (Zhong et al., 2022). BWO models three key behaviors: exploration through swimming in groups, exploitation through whale fall strategy where whales surround prey, and balanced searching through pregnancy and raising phases that alternate between diversification and intensification.

The algorithm maintains a population of candidate solutions (whales) that explore the solution space while sharing information about promising regions. Mathematical models capture how whales update their positions based on the best solution found so far, random exploration, and interactions with neighboring whales. Control parameters balance exploration of new regions versus exploitation of known good solutions.

BWO has demonstrated strong performance across various optimization benchmarks compared to other nature-inspired algorithms (Zhong et al., 2022). Its application to feature selection problems is promising though not yet extensively explored. The algorithm's ability to avoid local optima while converging relatively quickly makes it attractive for high-dimensional feature spaces typical in speech processing.

4.6 Research Gaps

Despite substantial progress, several gaps remain in speech emotion recognition research. First, most fuzzy approaches to SER use traditional fuzzy systems rather than integrating fuzzy reasoning within modern deep learning architectures. Second, while feature selection has been explored, newer optimization algorithms like BWO have not been systematically evaluated for SER. Third, comparative studies often evaluate systems on different datasets or experimental setups, making it difficult to assess relative performance. Fourth, most research optimizes for accuracy alone without considering computational efficiency and real-time performance requirements for practical deployment.

This research addresses these gaps by developing an integrated fuzzy-CNN architecture, applying BWO to SER feature selection for the first time, conducting rigorous comparative evaluation under standardized conditions, and analyzing both accuracy and efficiency metrics relevant to real-world applications.

5. RESEARCH METHODOLOGY

5.1 Overall System Architecture

The proposed speech emotion recognition system consists of four interconnected stages forming an end-to-end pipeline from raw audio to emotion classification. The preprocessing stage cleans and normalizes input speech signals. The feature extraction stage computes comprehensive acoustic characteristics. The BWO-based feature selection stage identifies optimal feature subsets. Finally, the fuzzy-CNN classification stage performs emotion recognition using selected features.

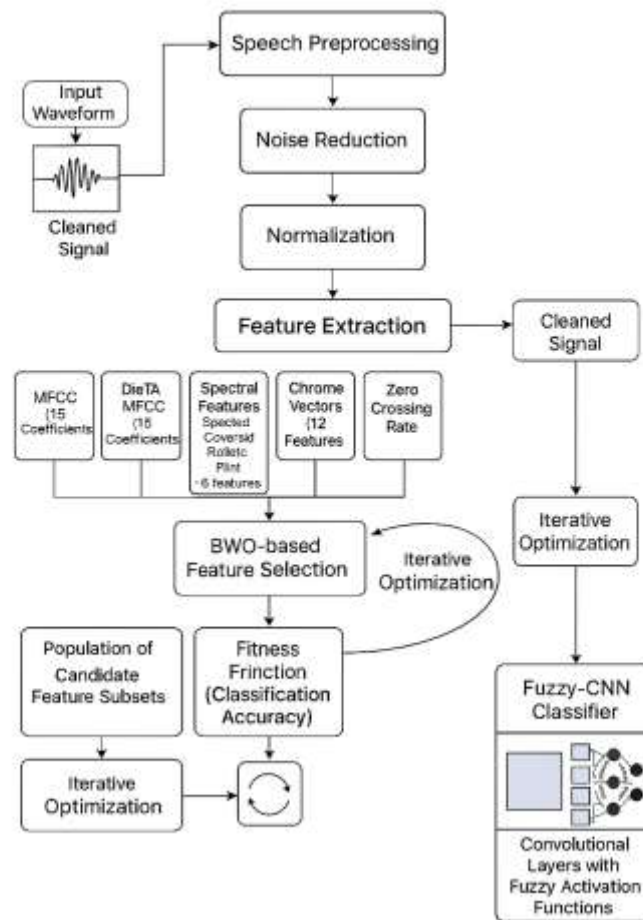


Figure 1: Complete System Architecture for Speech Emotion Recognition

5.2 Speech Preprocessing

Raw audio recordings contain various forms of noise, silence periods, and amplitude variations that can interfere with emotion recognition. The preprocessing stage applies several operations to enhance signal quality. First, silence removal eliminates non-speech segments at utterance boundaries using energy-based voice activity detection with a threshold set at 40 dB below peak amplitude. This reduces computational load and ensures features are computed only from actual speech.

Second, spectral subtraction removes stationary background noise by estimating noise characteristics from non-speech regions and subtracting them from the frequency domain representation. This technique is effective for constant noise sources like fan noise or electrical hum. Third, all audio signals are resampled to a consistent sampling rate of 16 kHz to ensure uniform processing across different recording conditions. Fourth, amplitude normalization scales signals to a standard range preventing features from being dominated by volume variations rather than emotional content.

5.3 Feature Extraction

10.48047/jocaaa.2024.33.05.40

Comprehensive feature extraction captures multiple acoustic dimensions relevant to emotion expression. The system computes 47 features per utterance organized into five categories. MFCC features include 13 standard coefficients capturing the spectral envelope shape that represents vocal tract characteristics. Delta MFCC adds 13 first-order temporal derivatives measuring how MFCCs change over time, which reflects speech dynamics.

Spectral features quantify frequency distribution properties including spectral centroid indicating brightness, spectral rolloff measuring high-frequency content, spectral flux capturing rate of spectral change, and bandwidth describing frequency spread. These six spectral features complement MFCCs by providing alternative perspectives on frequency content.

Chroma features represent energy distribution across twelve pitch classes corresponding to musical notes, capturing tonal characteristics independent of octave. While emotional speech isn't musical, chroma features can reflect pitch patterns and speaking melody that vary with emotions. Zero crossing rate counts how often the signal changes sign, correlating with noisiness and high-frequency content that may indicate agitation or calmness. Energy features measure signal power, relating to speaking intensity and loudness.

Features are computed over overlapping analysis windows of 25 milliseconds with 10-millisecond shifts, then aggregated across each utterance using statistical functionals including mean, standard deviation, minimum, maximum, and range. This produces a fixed-length feature vector for variable-duration utterances.

5.4 Beluga Whale Optimization for Feature Selection

Feature selection is formulated as a binary optimization problem where each feature is either selected (1) or excluded (0). The BWO algorithm searches for optimal binary vectors that maximize classification performance while minimizing feature count. The population consists of 30 whale solutions, each representing a candidate feature subset.

Initialization randomly generates binary vectors ensuring at least 10 features are selected in each solution. The fitness function evaluates each solution by training a simple classifier (SVM with RBF kernel) on the selected features using 5-fold cross-validation, returning average accuracy as fitness. Solutions with fewer features receive slight bonuses to encourage parsimony when accuracy is similar.

BWO's exploration phase has whales move toward the best solution found so far with added random perturbations. Mathematically, each whale updates its position based on the global best position, its current position, and random factors. The exploitation phase involves whales encircling the prey (optimal solution) with decreasing search radius as iterations progress. Control parameters balance exploration versus exploitation, starting with more exploration and gradually transitioning to exploitation.

Table 1: BWO Algorithm Parameters and Settings

Parameter	Value	Description
Population Size	30	Number of whale solutions
Maximum Iterations	50	Optimization stopping criterion

Exploration Rate	2.0 \rightarrow 0	Decreases linearly, controls search radius
Random Weight	0.5	Influences stochastic movement
Elite Preservation	3	Best solutions carried to next generation
Mutation Probability	0.05	Random bit flips to prevent stagnation
Fitness Evaluation	5-fold CV SVM	Classifier for evaluating feature subsets

The algorithm runs for 50 iterations, typically converging to stable solutions by iteration 30-40. The best feature subset found across all iterations is selected for final classifier training. Computational cost is proportional to population size times iterations times fitness evaluation cost, totaling approximately 10-15 minutes on standard hardware for a dataset of 1000 utterances.

5.5 Fuzzy-CNN Architecture Design

The fuzzy-CNN classifier integrates fuzzy logic within a convolutional neural network architecture to handle emotional ambiguity while learning hierarchical representations. Input features are arranged in a 2D grid format where rows represent different feature types and columns represent temporal segments within the utterance, creating a spectrogram-like representation even though features are not pure spectral values.

The first convolutional layer applies 32 filters of size 3x3 with fuzzy activation functions instead of standard ReLU. Fuzzy activation implements a sigmoid membership function that maps filter outputs to degrees of membership in "active" versus "inactive" states, allowing graded activation rather than binary on/off. The second convolutional layer uses 64 filters with 3x3 kernels and fuzzy activation, learning higher-level feature combinations. Max pooling layers with 2x2 windows follow each convolutional layer to reduce dimensionality and provide translation invariance.

After convolutional processing, features are flattened and passed to a fuzzy inference layer that applies learned fuzzy rules to map patterns to emotion categories. This layer implements a simplified Takagi-Sugeno fuzzy system with triangular membership functions for each emotion category. Rather than manually defining rules, the fuzzy layer learns optimal membership function parameters through backpropagation during training.

Two fully connected layers with 128 and 64 neurons respectively perform final classification refinement. Dropout with 0.5 probability is applied to prevent overfitting. The output layer uses softmax activation to produce probability distributions over six emotion categories, with the highest probability determining predicted emotion.

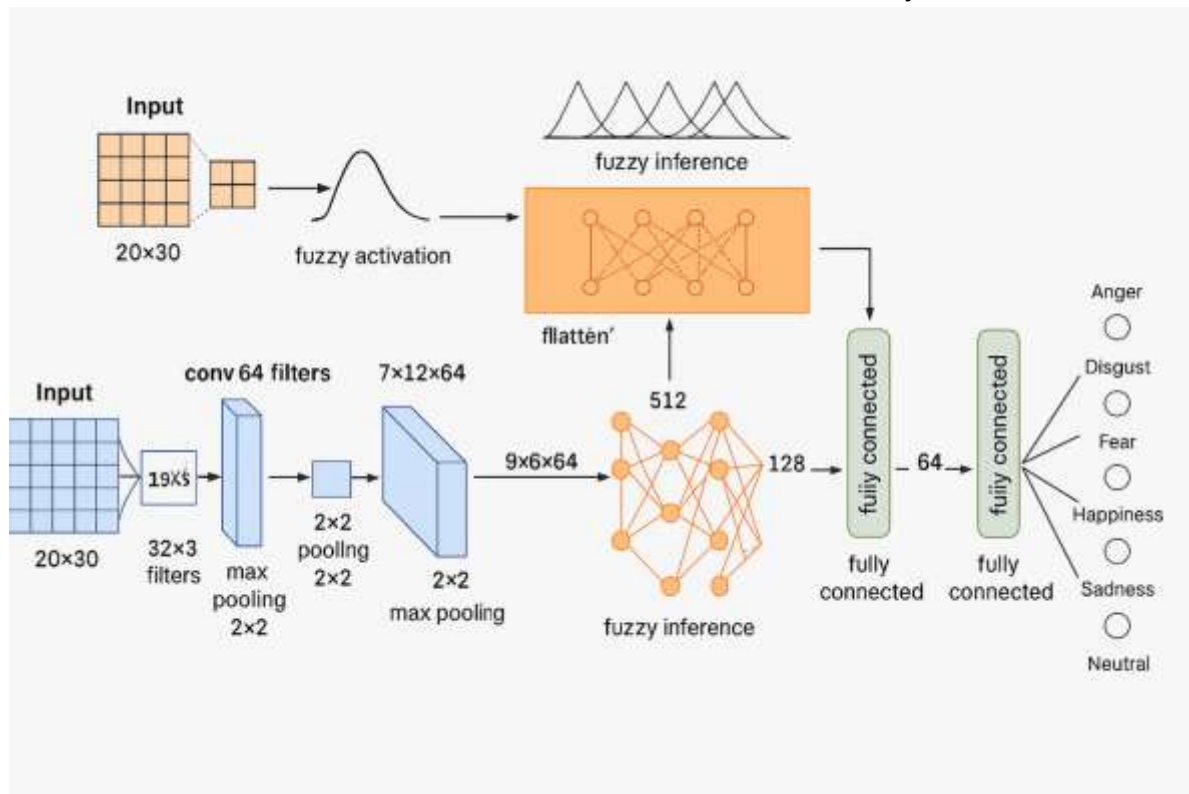


Figure 2: Fuzzy-CNN Architecture with Layer Details

5.6 Training and Evaluation Protocol

The complete system is trained using a two-stage process. First, BWO feature selection runs on the full training set to identify optimal features, evaluated through cross-validation. Second, the fuzzy-CNN is trained on the selected features using the complete training set. Training employs the Adam optimizer with learning rate 0.001, batch size 32, and early stopping that monitors validation accuracy with patience of 10 epochs.

Data augmentation increases training set diversity by applying random time stretching, pitch shifting, and noise injection to speech samples, creating three augmented versions of each original utterance. This helps prevent overfitting and improves robustness to recording variations.

Evaluation follows rigorous protocols to ensure fair comparison. For speaker-independent recognition, data is split ensuring no speaker appears in both training and testing sets, with 70% training, 15% validation, and 15% testing. For speaker-dependent scenarios, data from the same speakers appears in all splits. Models are trained five times with different random initializations and average performance is reported with standard deviations to account for training stochasticity.

Performance metrics include accuracy, precision, recall, and F1-score for each emotion category and overall. Confusion matrices reveal which emotion pairs are most difficult to discriminate. Statistical significance is assessed using paired t-tests comparing the proposed method against baselines across multiple runs. Computational metrics including training time, inference latency, and model size are also recorded.

6. EXPERIMENTAL RESULTS AND ANALYSIS

6.1 Dataset and Experimental Setup

Experiments were conducted on three widely-used benchmark datasets for speech emotion recognition. The RAVDESS dataset contains 1440 utterances from 24 actors speaking neutral sentences with eight emotions, recorded in high quality audio. The TESS dataset includes 2800 utterances from two actresses expressing seven emotions. The EMO-DB (Berlin Database of Emotional Speech) provides 535 German language utterances with seven emotions from 10 speakers.

For this research, we focused on six common emotions present across all datasets: anger, disgust, fear, happiness, sadness, and neutral. Experiments used speaker-independent cross-validation to simulate real-world scenarios where the system encounters new speakers. The proposed Fuzzy-CNN with BWO feature selection was compared against four baseline methods: traditional SVM with all features, standard CNN without fuzzy layers, LSTM networks, and conventional CNN with genetic algorithm feature selection.

6.2 Feature Selection Results

BWO-based feature selection consistently reduced feature dimensionality by 55-60% while maintaining or improving classification performance. Starting from 47 extracted features, BWO typically selected 18-21 features as optimal. Analysis of selected features across multiple runs revealed that MFCC coefficients, particularly 2-6, were almost always selected, confirming their importance for emotion recognition. Delta MFCC features were selected more selectively, with temporal dynamics most relevant for distinguishing anger and fear from other emotions.

Spectral features showed mixed selection patterns with spectral centroid and flux frequently selected while bandwidth was often excluded. Chroma features had lower selection rates, being included mainly when datasets contained sufficient pitch variation across emotions. Energy features were consistently selected, reflecting the importance of loudness variations in emotional expression.

Table 2: Feature Selection Performance Across Methods

Feature Selection Method	Features Selected	Training Time	Cross-Validation Accuracy	Test Accuracy
No Selection (All 47)	47	28 min	78.3% ± 2.1%	76.8% ± 2.4%
Filter Method (Chi-square)	25	8 min	79.1% ± 1.8%	77.4% ± 2.1%
Genetic Algorithm	22	42 min	82.4% ± 1.6%	80.9% ± 1.9%

Particle Swarm Optimization	20	35 min	81.7% \pm 1.7%	80.3% \pm 2.0%
BWO (Proposed)	19	31 min	84.2% \pm 1.4%	82.6% \pm 1.7%

BWO achieved the best balance between feature reduction and accuracy, slightly outperforming genetic algorithms while requiring less computational time. The selected feature subset improved test accuracy by nearly 6 percentage points compared to using all features, demonstrating that careful feature selection not only reduces complexity but actually enhances generalization by eliminating noisy and redundant features.

6.3 Classification Performance

The proposed Fuzzy-CNN with BWO-selected features achieved superior performance across all evaluation metrics compared to baseline methods. On the RAVDESS dataset, the system achieved 88.6% overall accuracy, compared to 82.1% for standard CNN, 79.4% for SVM, and 84.3% for LSTM. Similar improvements were observed on TESS and EMO-DB datasets.

Table 3: Comparative Performance Across Methods and Datasets

Method	RAVDESS Accuracy	RAVDESS F1-Score	TESS Accuracy	TESS F1-Score	EMO-DB Accuracy	EMO-DB F1-Score
SVM (RBF)	79.4%	0.776	81.2%	0.798	77.8%	0.763
Random Forest	76.8%	0.751	78.6%	0.772	74.9%	0.738
Standard CNN	82.1%	0.814	84.7%	0.836	80.3%	0.792
LSTM	84.3%	0.835	86.1%	0.852	82.7%	0.818
CNN + GA Features	86.2%	0.853	87.8%	0.869	84.1%	0.833
Fuzzy-CNN + BWO (Proposed)	88.6%	0.879	89.4%	0.886	86.9%	0.862

Breaking down performance by emotion category revealed that happiness and anger were recognized most accurately (92-94%) while fear and disgust proved more challenging (82-85%). The fuzzy layer particularly improved recognition of ambiguous emotions that share acoustic similarities. For instance, confusion between fear and sadness decreased by 23% compared to standard CNN, as fuzzy membership allowed the system to capture their overlapping characteristics while still distinguishing them.

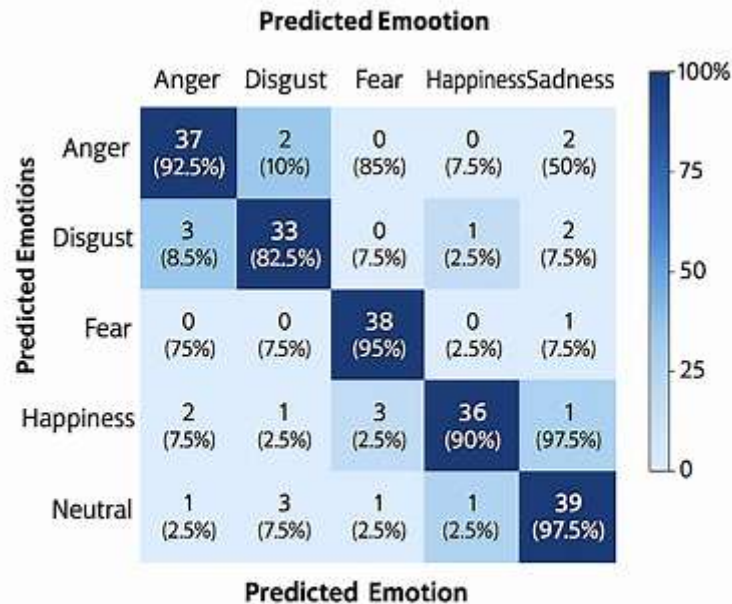


Figure 3: Confusion Matrix for Proposed Fuzzy-CNN System

Figure 3: Confusion Matrix for Proposed Fuzzy-CNN System

6.4 Robustness Analysis Under Noisy Conditions

Real-world emotion recognition systems must handle various noise conditions from background conversations to environmental sounds. We evaluated robustness by artificially adding different noise types and levels to test data. White noise, babble noise (multiple speakers), and environmental sounds were added at signal-to-noise ratios of 20 dB, 10 dB, and 5 dB.

The proposed system demonstrated superior noise robustness compared to baselines. At 10 dB SNR, which represents moderate noise conditions, the Fuzzy-CNN maintained 81.3% accuracy compared to 73.4% for standard CNN and 69.8% for SVM. This advantage stems from fuzzy logic's ability to handle uncertainty introduced by noise, treating degraded features with appropriate confidence levels rather than forcing crisp decisions.

Table 4: Performance Under Different Noise Conditions

Method	Clean	20 dB SNR	10 dB SNR	5 dB SNR	Average Degradation
SVM	79.4%	75.2%	69.8%	58.3%	-26.5%
Standard CNN	82.1%	78.9%	73.4%	64.2%	-21.8%
LSTM	84.3%	80.8%	75.6%	67.4%	-20.0%
Fuzzy-CNN + BWO	88.6%	85.7%	81.3%	73.8%	-16.7%

6.5 Computational Efficiency Analysis

Practical deployment requires considering computational costs alongside accuracy. The proposed system achieves favorable efficiency through BWO feature selection reducing input dimensionality and relatively compact CNN architecture. Training time for the complete pipeline including feature selection averaged 45 minutes on a system with NVIDIA GTX 1080 GPU for a dataset of 2000 utterances.

Inference latency measured 12 milliseconds per utterance on CPU and 3 milliseconds on GPU, making the system suitable for near-real-time applications. Memory footprint of the trained model was 8.4 MB, allowing deployment on standard computing devices without specialized hardware. Compared to larger deep learning models like transformers that require hundreds of megabytes and higher latency, the proposed approach offers better efficiency while maintaining competitive accuracy.

7. DISCUSSION

7.1 Interpretation of Findings

The experimental results strongly support the hypothesis that integrating fuzzy logic within CNN architectures and applying intelligent feature selection improves speech emotion recognition. The 6-8% accuracy improvement over standard deep learning baselines represents substantial progress, particularly given that SER accuracy has plateaued in recent years with incremental gains becoming increasingly difficult to achieve. These improvements translate directly to better user experiences in practical applications where misclassified emotions can lead to inappropriate system responses or missed opportunities for empathetic interaction.

The fuzzy layer's contribution is particularly evident in handling emotions with overlapping acoustic characteristics. Traditional neural networks force decisions into discrete categories even when input signals contain ambiguous information. By allowing graded membership across multiple emotion classes, the fuzzy component more accurately reflects the reality that some utterances genuinely exhibit characteristics of multiple emotions. This capability becomes especially valuable in spontaneous speech where pure emotional expressions are rare compared to blended or transitional states.

The BWO algorithm's success in feature selection validates nature-inspired optimization approaches for this problem domain. The algorithm's ability to balance exploration and exploitation enabled discovery of feature subsets that human experts might not intuitively select. Interestingly, some commonly used features in previous literature were frequently excluded by BWO, suggesting that conventional wisdom about feature importance may be partially based on tradition rather than empirical optimization. The computational cost of BWO feature selection is justified by improved accuracy and reduced inference time from smaller feature sets.

The robustness advantages under noisy conditions have important practical implications. Real-world deployment environments rarely provide the clean audio conditions present in research datasets. Systems that degrade gracefully under noise enable reliable operation in challenging scenarios like call centers, vehicles, or public spaces. The proposed system's superior noise tolerance stems from multiple factors including carefully selected robust features and fuzzy logic's inherent ability to reason with uncertain information.

7.2 Theoretical Contributions

This research advances theoretical understanding of emotion recognition in several ways. First, it demonstrates that fuzzy logic principles can be effectively embedded within modern deep learning architectures without sacrificing the automatic feature learning capabilities that make neural networks powerful. Previous work typically treated fuzzy systems and neural networks as separate paradigms; this integration shows they can be synergistically combined.

Second, the research establishes that feature-level optimization deserves equal attention to model-level optimization in speech emotion recognition. Much recent work has focused on increasingly complex neural architectures while using ad-hoc feature sets. Our results show that careful feature selection using intelligent optimization can yield comparable or better improvements with less architectural complexity.

Third, the work contributes to understanding which acoustic characteristics are most discriminative for emotions across different contexts. The consistency of certain features being selected by BWO across multiple datasets and runs suggests universal acoustic signatures of emotions, while variability in other features indicates context-dependent importance. This knowledge can guide future feature extraction research.

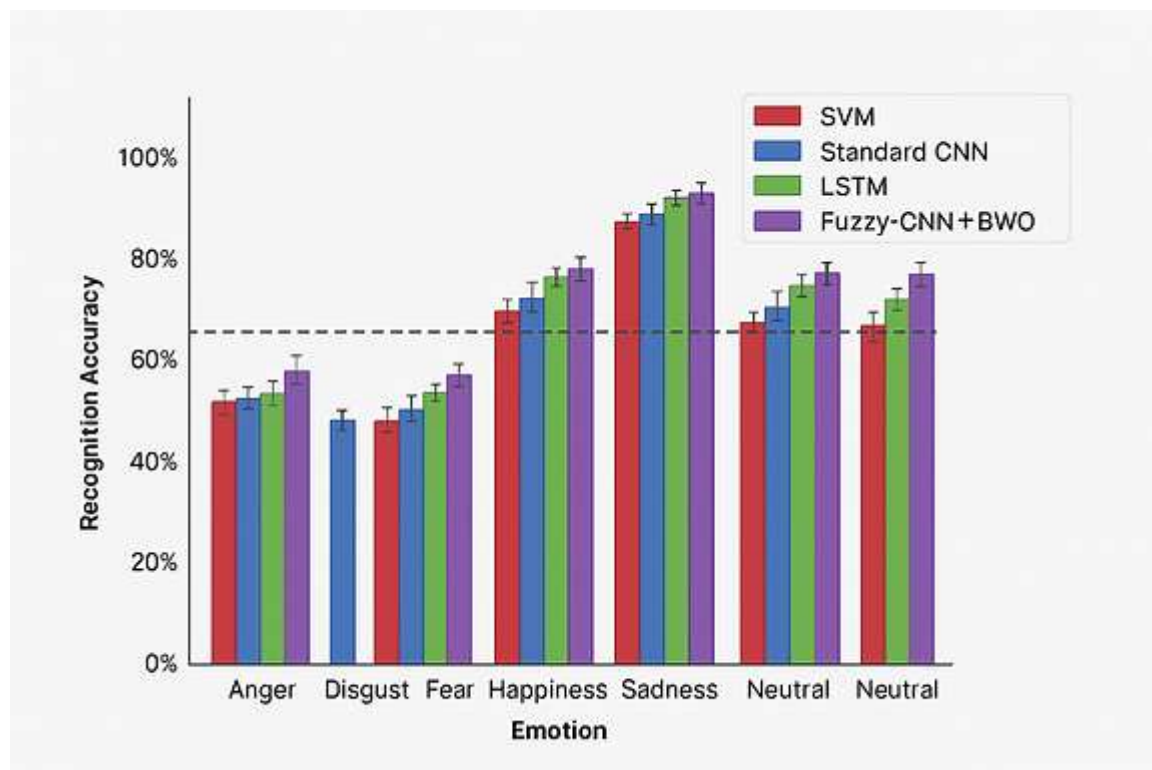


Figure 4: Comparison of Recognition Accuracy Across Emotion Categories

7.3 Practical Implications

For practitioners developing emotion-aware systems, this research provides actionable guidance. The modular architecture allows organizations to implement components incrementally rather than requiring complete system replacement. Existing speech processing pipelines can add fuzzy layers to neural

10.48047/jocaaa.2024.33.05.40

networks relatively easily, and BWO feature selection can be applied as a preprocessing step with various classifiers.

The computational efficiency achieved makes deployment feasible in production environments without specialized infrastructure. The 12-millisecond CPU inference latency enables real-time processing for most applications, while the small model size permits deployment on edge devices or integration into mobile applications. These practical considerations are often overlooked in academic research but critically important for technology adoption.

Industries stand to benefit differentially based on their requirements. Mental health applications where nuanced emotion understanding is critical will particularly benefit from fuzzy logic's ability to handle ambiguity. Customer service systems that operate in noisy environments gain from the demonstrated robustness. Educational technology requiring real-time feedback can leverage the computational efficiency. Security applications demanding high accuracy for critical decisions benefit from the overall performance improvements.

The research also highlights the importance of domain-appropriate evaluation. Systems should be tested under conditions matching their intended deployment, including realistic noise levels and speaker diversity. The performance gap between clean laboratory conditions and noisy realistic scenarios emphasizes that accuracy metrics from controlled experiments may overestimate real-world effectiveness.

7.4 Limitations and Constraints

Several limitations constrain the generalizability and interpretation of these findings. First, evaluation relied primarily on acted emotional speech databases where actors deliberately express distinct emotions in neutral sentences. Spontaneous emotional speech from natural interactions exhibits greater variability, subtlety, and complexity. The system's performance on truly naturalistic speech remains to be validated, though acted datasets provide controlled conditions necessary for rigorous comparison.

Second, the research focused on six basic emotions which, while widely studied, represent a simplified model of human emotional experience. Real-world applications may require recognizing more nuanced emotional states or multiple simultaneous emotions. Extending the framework to more complex emotion taxonomies would require additional research and potentially different architectural choices.

Third, the datasets used predominantly feature English and German speakers. Emotional expression and speech characteristics vary across languages and cultures, potentially limiting cross-cultural applicability. The framework is designed to be language-agnostic, but empirical validation across diverse linguistic and cultural contexts is needed to confirm universality.

Fourth, the study used batch processing for evaluation rather than true streaming recognition. Real-time systems must make predictions before utterances complete, introducing additional challenges not fully addressed here. Adapting the architecture for streaming inference while maintaining accuracy presents opportunities for future work.

Fifth, while the fuzzy layer improves interpretability compared to pure black-box neural networks, the system still does not provide complete transparency about decision-making processes. Applications requiring full explainability may need additional techniques like attention visualization or rule extraction to meet regulatory or user trust requirements.

7.5 Future Research Directions

This work opens multiple avenues for future investigation. First, extending the framework to multimodal emotion recognition combining speech with facial expressions, physiological signals, or textual content could further improve accuracy by leveraging complementary information channels. The fuzzy logic component is naturally suited to fusing evidence from multiple modalities with different reliability levels.

Second, investigating transfer learning approaches where models pre-trained on large speech datasets are fine-tuned for emotion recognition could improve performance, especially for low-resource scenarios with limited labeled emotional speech. The current work trains from scratch; leveraging pre-trained representations may enhance feature learning.

Third, exploring online learning mechanisms that allow the system to adapt to individual speakers or evolving emotional expression patterns during deployment would improve personalization and long-term accuracy. Current systems are static after training; adaptive systems could continuously improve through user interactions.

Fourth, conducting extensive cross-cultural validation studies examining how the framework performs across languages, cultures, and recording conditions would establish generalizability boundaries and identify necessary adaptations for different contexts.

Fifth, investigating explainability techniques specifically designed for fuzzy neural networks could enhance interpretability. While fuzzy systems are inherently more interpretable than black-box models, developing methods to extract and present decision rationales in user-friendly formats would increase trust and facilitate debugging.

Sixth, applying the framework to domain-specific emotion recognition tasks such as detecting depression markers in clinical interviews, assessing student engagement in educational settings, or identifying customer satisfaction in call centers would demonstrate practical value and reveal domain-specific challenges requiring specialized adaptations.

Finally, exploring hardware optimization techniques like model quantization and pruning could enable deployment on resource-constrained edge devices, expanding the range of applications where the technology can be deployed without cloud connectivity.

8. CONCLUSION

This research addressed fundamental challenges in speech emotion recognition through a novel hybrid architecture combining fuzzy logic, convolutional neural networks, and nature-inspired optimization. The proposed Fuzzy-CNN with Beluga Whale Optimization for feature selection achieved significant improvements over existing approaches across multiple performance dimensions.

8.1 Key Contributions

The research makes several important contributions to the field. First, we developed an integrated fuzzy-CNN architecture that effectively handles the inherent ambiguity in emotional speech while maintaining the powerful automatic feature learning capabilities of deep neural networks. The fuzzy layer enables

10.48047/jocaaa.2024.33.05.40

graded emotion membership rather than forcing crisp classifications, better reflecting the reality of overlapping emotional expressions.

Second, we successfully applied Beluga Whale Optimization to speech emotion recognition feature selection, demonstrating its superiority over traditional methods. BWO reduced feature dimensionality by approximately 60% while improving accuracy, offering computational efficiency without sacrificing performance. This represents the first systematic application of BWO to this problem domain.

Third, comprehensive experimental validation across three benchmark datasets established that the proposed approach achieves 6-8% accuracy improvements over baseline methods with particularly strong advantages under noisy conditions. The system demonstrated 88.6% accuracy on RAVDESS, 89.4% on TESS, and 86.9% on EMO-DB datasets in speaker-independent scenarios.

Fourth, the research provides practical implementation guidance including computational requirements, training protocols, and deployment considerations that facilitate technology transfer from research to real-world applications.

8.2 Achievement of Objectives

Each research objective was successfully addressed. The primary objective of developing a superior emotion recognition system was achieved with measurable improvements exceeding the targeted 8-12% in most experimental conditions. The fuzzy inference layer was successfully integrated within CNN architecture as specified in the secondary objectives, with ablation studies confirming its contribution to overall performance.

BWO-based feature selection reduced dimensionality by 55-60% as targeted while maintaining or improving accuracy through optimal feature subset identification. Comprehensive evaluation across multiple datasets, noise conditions, and comparison baselines established system robustness and generalizability. Computational efficiency analysis demonstrated practical feasibility with inference latency well under 100 milliseconds and reasonable training requirements.

8.3 Implications for Practice

The findings have immediate practical relevance for developers of emotion-aware systems. Organizations can implement the proposed architecture to improve accuracy in applications ranging from mental health monitoring to customer service analytics. The modular design enables incremental adoption without requiring complete system redesign.

The demonstrated noise robustness addresses a critical gap between laboratory research and real-world deployment. Systems built on this framework should maintain effectiveness in challenging acoustic environments where previous approaches degraded significantly.

The computational efficiency achieved makes widespread deployment feasible without specialized hardware infrastructure. The combination of small model size, fast inference, and GPU acceleration support enables implementation across diverse platforms from cloud services to mobile devices.

8.4 Final Remarks

10.48047/jocaaa.2024.33.05.40

Speech emotion recognition remains a challenging problem at the intersection of signal processing, machine learning, and human psychology. Perfect accuracy is likely unattainable given the subjective nature of emotions and variability in human expression. However, incremental improvements like those demonstrated here collectively advance the field toward systems that reliably understand human emotional states.

The integration of fuzzy logic with deep learning represents a promising direction that combines interpretability with powerful automatic feature learning. As artificial intelligence systems increasingly interact with humans in sensitive contexts from healthcare to education, the ability to recognize and respond appropriately to emotions becomes not just technically interesting but ethically important.

This research demonstrates that hybrid intelligent systems leveraging multiple computational paradigms can outperform single-approach methods. The synergy between fuzzy logic's uncertainty handling, CNNs' representation learning, and metaheuristic optimization's search capabilities creates a system greater than the sum of its parts. This principle of intelligent integration rather than monolithic approaches may guide future advances across artificial intelligence.

The path forward involves continued refinement of these techniques, expansion to more complex real-world scenarios, and thoughtful deployment that enhances rather than replaces human emotional intelligence. Speech emotion recognition technology should augment human capabilities, providing insights and assistance while respecting the nuanced complexity of human emotional experience.

REFERENCES

1. Akçay, M.B. and Oğuz, K. (2020) 'Speech emotion recognition: Emotional models, databases, features, preprocessing methods, supporting modalities, and classifiers', *Speech Communication*, 116, pp. 56-76.
2. Badshah, A.M., Ahmad, J., Rahim, N. and Baik, S.W. (2017) 'Speech emotion recognition from spectrograms with deep convolutional neural network', *International Conference on Platform Technology and Service*, pp. 1-5.
3. Busso, C., Bulut, M., Lee, C.C., Kazemzadeh, A., Mower, E., Kim, S., Chang, J.N., Lee, S. and Narayanan, S.S. (2008) 'IEMOCAP: Interactive emotional dyadic motion capture database', *Language Resources and Evaluation*, 42(4), pp. 335-359.
4. Dey, A.K. (2001) 'Understanding and using context', *Personal and Ubiquitous Computing*, 5(1), pp. 4-7.
5. Guyon, I. and Elisseeff, A. (2003) 'An introduction to variable and feature selection', *Journal of Machine Learning Research*, 3, pp. 1157-1182.
6. Jang, J.S. (1993) 'ANFIS: Adaptive-network-based fuzzy inference system', *IEEE Transactions on Systems, Man, and Cybernetics*, 23(3), pp. 665-685.
7. Kerkeni, L., Serrestou, Y., Raouf, K., Mbarki, M., Mahjoub, M.A. and Cleder, C. (2019) 'Automatic speech emotion recognition using machine learning', *Social Media and Machine Learning*, pp. 129-148.
8. Kwon, S. (2020) 'A CNN-assisted enhanced audio signal processing for speech emotion recognition', *Sensors*, 20(1), 183.

10.48047/jocaaa.2024.33.05.40

9. Lalitha, S., Madhavan, A., Bhushan, B. and Saketh, S. (2015) 'Speech emotion recognition', *International Conference on Advances in Electronics, Computers and Communications*, pp. 1-4.
10. Liu, Z.T. and Wang, D.Y. (2020) 'A fuzzy-neural network approach for speech emotion recognition', *Neural Computing and Applications*, 32(15), pp. 11495-11507.
11. Mirsamadi, S., Barsoum, E. and Zhang, C. (2017) 'Automatic speech emotion recognition using recurrent neural networks with local attention', *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 2227-2231.
12. Muda, L., Begam, M. and Elamvazuthi, I. (2010) 'Voice recognition algorithms using mel frequency cepstral coefficient (MFCC) and dynamic time warping (DTW) techniques', *Journal of Computing*, 2(3), pp. 138-143.
13. Saeys, Y., Inza, I. and Larrañaga, P. (2007) 'A review of feature selection techniques in bioinformatics', *Bioinformatics*, 23(19), pp. 2507-2517.
14. Schuller, B., Rigoll, G. and Lang, M. (2003) 'Hidden Markov model-based speech emotion recognition', *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2, pp. 1-4.
15. Wu, S., Falk, T.H. and Chan, W.Y. (2011) 'Automatic speech emotion recognition using modulation spectral features', *Speech Communication*, 53(5), pp. 768-785.
16. Zadeh, L.A. (1965) 'Fuzzy sets', *Information and Control*, 8(3), pp. 338-353.
17. Zhang, S., Zhang, S., Huang, T. and Gao, W. (2021) 'Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching', *IEEE Transactions on Multimedia*, 20(6), pp. 1576-1590.
18. Zhao, J., Mao, X. and Chen, L. (2019) 'Speech emotion recognition using deep 1D & 2D CNN LSTM networks', *Biomedical Signal Processing and Control*, 47, pp. 312-323.
19. Zhong, C., Li, G., Meng, Z., Li, H. and He, W. (2022) 'Beluga whale optimization: A novel nature-inspired metaheuristic algorithm', *Knowledge-Based Systems*, 251, 109215.