

Binary Speech Emotion Recognition: A Comparative Study of SVM, CNN, and Bi-LSTM Approaches

Progress Report

Ammar Qurthuby*, Habibi*

*Department of Informatics Engineering, Universitas Syiah Kuala

Email: {ammar22, habibi123}@mhs.usk.ac.id

Abstract—This progress report presents the implementation status and preliminary results of our speech emotion recognition research project. We have successfully merged four publicly available databases (CREMA-D, RAVDESS, SAVEE, and TESS) containing 12,162 audio samples and implemented a binary emotion classification system. The dataset has been preprocessed and split into Negative (angry, disgust, fear, sad) comprising 7,692 samples (63.2%) and Non-Negative (happy, neutral, surprise) comprising 4,470 samples (36.8%). We have completed the implementation of dual feature extraction strategies: 80-dimensional statistical features for traditional machine learning and 100x40 sequential MFCC features for deep learning models. The SVM baseline model has been trained and evaluated, achieving 79.57% test accuracy with balanced performance across both classes. Per-class results show Precision of 0.82/0.75, Recall of 0.86/0.68, and F1-Score of 0.84/0.71 for Negative/Non-Negative classes respectively. Deep learning models (CNN and Bi-LSTM) are currently under development. This report details the completed work, preliminary findings, and remaining tasks toward project completion.

Index Terms—Speech Emotion Recognition, Binary Classification, Support Vector Machine, Progress Report, MFCC Features, Dataset Merging

I. INTRODUCTION

A. Project Overview

This progress report documents the implementation of a comparative study on binary speech emotion recognition using three machine learning approaches: Support Vector Machine (SVM), 1D Convolutional Neural Network (CNN), and Bidirectional Long Short-Term Memory (Bi-LSTM). As proposed in our initial plan, we aim to classify emotions into two categories: Negative (angry, disgust, fear, sad) and Non-Negative (happy, neutral, surprise) using a merged dataset from multiple sources.

B. Objectives Recap

The primary objectives of this research are:

- 1) Develop a comprehensive merged dataset from four publicly available databases
- 2) Implement and evaluate traditional machine learning (SVM) as baseline
- 3) Develop deep learning architectures (CNN and Bi-LSTM) for comparison
- 4) Analyze performance trade-offs between different approaches
- 5) Provide deployment recommendations for practical applications

C. Progress Summary

As of Week 6 (December 2025), we have successfully completed several major milestones:

- **Dataset Collection and Preprocessing:** Merged four databases and preprocessed 12,162 audio samples with stratified train-test split (100% complete)
- **Feature Extraction Pipeline:** Implemented dual feature strategies for traditional ML and deep learning approaches (100% complete)
- **SVM Baseline Development:** Trained and evaluated SVM with RBF kernel, achieving 79.57% test accuracy (100% complete)
- **CNN Architecture Design:** Designed and partially implemented 1D-CNN model with three convolutional blocks (80% complete)
- **Bi-LSTM Architecture Planning:** Completed architecture design with bidirectional LSTM layers (60% complete)

The following sections detail our methodology, preliminary experimental results, and remaining work.

II. LITERATURE REVIEW

A. Speech Emotion Recognition Approaches

Speech emotion recognition remains an active research area with applications in human-computer interaction and mental health monitoring [?]. El Ayadi et al. [?] provided a comprehensive survey demonstrating that effective SER systems require careful consideration of feature extraction, classification algorithms, and dataset characteristics.

Traditional approaches using Support Vector Machines have shown consistent performance across various emotion recognition tasks [?]. More recent deep learning approaches, particularly CNNs operating on spectrograms [?], [?] and RNNs capturing temporal dependencies [?], [?], have achieved state-of-the-art results on several benchmarks.

B. Binary vs. Multi-Class Classification

While most SER research focuses on multi-class emotion recognition (4-8 categories), binary classification into positive/negative valence has gained attention for clinical applications [?]. Binary classification typically achieves higher accuracy than fine-grained emotion recognition and provides sufficient information for mental health screening applications [?].

C. Feature Extraction Methods

Mel-Frequency Cepstral Coefficients (MFCC) remain the most widely adopted features for SER due to their effectiveness in representing spectral envelope information relevant to human auditory perception [?]. The choice between statistical aggregation (mean, standard deviation) and preserving temporal sequences depends on the classification algorithm: traditional classifiers require fixed-length vectors, while deep learning models can process variable-length sequences [?].

Our dual feature extraction strategy addresses both requirements, enabling fair comparison between traditional and deep learning approaches.

III. METHODOLOGY

A. Dataset Construction

1) *Source Databases*: We merged four publicly available speech emotion databases:

- 1) **CREMA-D** (Crowd-sourced Emotional Multimodal Actors Dataset) [?]: Contains 7,442 audio clips from 91 actors expressing six emotions with multiple intensity levels.
- 2) **RAVDESS** (Ryerson Audio-Visual Database of Emotional Speech and Song) [?]: Includes 1,440 speech files from 24 professional actors (12 male, 12 female) aged 21-33 years.
- 3) **SAVEE** (Surrey Audio-Visual Expressed Emotion Database) [?]: Comprises 480 utterances from 4 male actors in 7 emotion categories.
- 4) **TESS** (Toronto Emotional Speech Set) [?]: Provides 2,800 audio files from 2 female actors (aged 26 and 64) reading 200 target words in 7 emotions.

2) *Emotion Mapping and Distribution*: Original emotion labels from each database were mapped to our binary classification scheme:

- **Negative Class**: angry, disgust, fear, sad
- **Non-Negative Class**: happy, neutral, surprise

Table ?? presents the final dataset composition after merging and preprocessing.

TABLE I
MERGED DATASET DISTRIBUTION (12,162 SAMPLES)

Category	Emotion	Count	% Total
Negative	Angry	1,923	15.8%
	Disgust	1,923	15.8%
	Fear	1,923	15.8%
	Sad	1,923	15.8%
<i>Subtotal</i>		7,692	63.2%
Non-Negative	Happy	1,923	15.8%
	Neutral	1,895	15.6%
	Surprise	652	5.4%
	<i>Subtotal</i>	4,470	36.8%
Total		12,162	100%

The dataset exhibits moderate class imbalance with a 1.72:1 ratio (Negative:Non-Negative). Additionally, the surprise emo-

tion is notably underrepresented (652 samples, 5.4%), which may impact model performance on this specific emotion.

3) *Data Splitting*: The merged dataset was split using stratified sampling to maintain class balance:

- **Training set**: 9,730 samples (80% of total dataset)
- **Test set**: 2,432 samples (20% of total dataset)

Stratification ensures proportional representation of both emotion classes (Negative and Non-Negative) and all seven individual emotions in both training and test sets, enabling reliable and unbiased performance evaluation.

B. Feature Extraction

We implemented two feature extraction strategies to accommodate different model requirements:

1) *Statistical Features for SVM*: For traditional machine learning, we extracted 80-dimensional statistical features from each audio sample:

- **MFCC Statistics (40D)**: Mean and standard deviation of 20 MFCC coefficients computed over the entire utterance, capturing spectral envelope characteristics.
- **Acoustic Features (40D)**: Mean and standard deviation of five key acoustic properties:
 - Zero-crossing rate (2D) – indicates frequency content
 - Spectral centroid (2D) – represents brightness of sound
 - Spectral bandwidth (2D) – measures frequency range
 - Spectral rolloff (2D) – indicates high-frequency content
 - RMS energy (2D) – captures overall loudness
 - 15 additional spectral features (30D) – complementary spectral descriptors

All features were extracted using the `librosa` Python library (version 0.9.2) with consistent parameters: sampling rate of 22,050 Hz, frame length of 2,048 samples, and hop length of 512 samples.

2) *Sequential Features for Deep Learning*: For CNN and Bi-LSTM models, we extracted sequential MFCC representations that preserve temporal dynamics:

- **Frame-level Features**: 40 MFCC coefficients extracted per frame
- **Sequence Length**: Fixed to 100 time steps for uniform input shape
- **Padding Strategy**: Zero-padding applied to shorter sequences
- **Truncation**: Longer sequences trimmed to 100 frames from the center
- **Resulting Shape**: 100×40 matrices per audio sample
- **Normalization**: Z-score standardization applied: $(x - \mu)/\sigma$

This sequential representation allows deep learning models to learn temporal patterns and dependencies that are lost when using statistical aggregation (mean, standard deviation) required by traditional classifiers.

C. SVM Implementation (Completed)

1) *Architecture and Hyperparameters:* We implemented an SVM classifier with the following configuration:

- **Kernel:** Radial Basis Function (RBF)
- **Hyperparameter Optimization:** 5-fold cross-validation grid search
- **Search Space:**
 - Regularization parameter C : [0.1, 1, 10, 100]
 - Kernel coefficient γ : [0.001, 0.01, 0.1, 1]
- **Class Weights:** Balanced to handle class imbalance
- **Optimal Parameters:** $C = 10, \gamma = 0.01$

2) *Training Process:* The SVM classifier was trained on 80-dimensional statistical feature vectors extracted from 9,730 training samples. Grid search with 5-fold cross-validation explored 16 hyperparameter combinations (4 values of C \times 4 values of γ), requiring approximately 45 minutes of computation on our system (Intel Core i7 processor, 16GB RAM, no GPU acceleration). The optimal hyperparameters were selected based on maximizing cross-validation weighted F1-score to account for class imbalance. The final model with $C = 10$ and $\gamma = 0.01$ was then retrained on the entire training set for evaluation on the held-out test set.

D. CNN Architecture (In Progress)

We designed a 1D Convolutional Neural Network architecture for sequential MFCC processing:

- **Input Layer:** Accepts 100x40 MFCC sequences
- **Conv Block 1:** 64 filters with kernel size 5, ReLU activation, BatchNorm, MaxPooling (pool size 2)
- **Conv Block 2:** 128 filters with kernel size 5, ReLU activation, BatchNorm, MaxPooling (pool size 2)
- **Conv Block 3:** 256 filters with kernel size 3, ReLU activation, BatchNorm, MaxPooling (pool size 2)
- **Global Average Pooling:** Reduces spatial dimensions while retaining learned features
- **Dense Layer 1:** 256 units with ReLU activation and Dropout (rate: 0.5)
- **Output Layer:** 2 units with softmax activation for binary classification

Each convolutional block includes batch normalization for training stability and max pooling for dimensionality reduction. The model is currently being trained with Adam optimizer (learning rate: 0.001, $\beta_1 = 0.9, \beta_2 = 0.999$) for up to 100 epochs with early stopping (patience: 15 epochs monitoring validation loss).

E. Bi-LSTM Architecture (Planned)

The Bidirectional LSTM architecture has been designed to capture temporal dependencies in both forward and backward directions:

- **Input Layer:** Accepts 100x40 MFCC sequences
- **Bi-LSTM Layer 1:** 128 units (256 total: 128 forward + 128 backward), returns sequences
- **Bi-LSTM Layer 2:** 64 units (128 total: 64 forward + 64 backward), returns sequences

- **Bi-LSTM Layer 3:** 32 units (64 total: 32 forward + 32 backward)
- **Dropout:** Rate of 0.3 applied between LSTM layers to prevent overfitting
- **Dense Layer 1:** 64 units with ReLU activation and Dropout (rate: 0.4)
- **Output Layer:** 2 units with softmax activation for binary classification

The architecture will be implemented using TensorFlow/Keras with Adam optimizer (learning rate: 0.0001) and trained with early stopping (patience: 20 epochs). Implementation is scheduled for Week 7, following CNN model completion and evaluation.

IV. PRELIMINARY EXPERIMENTAL RESULTS

A. SVM Baseline Performance

Table ?? presents the classification performance of the optimized SVM model on the test set (2,432 samples).

TABLE II
SVM CLASSIFICATION RESULTS

Metric	Negative	Non-Negative	Avg.
Precision	0.82	0.75	0.80
Recall	0.86	0.68	0.80
F1-Score	0.84	0.71	0.80
Test Accuracy	79.57%		

B. Performance Analysis

The SVM baseline achieved 79.57% test accuracy, which represents a strong foundation for comparison with upcoming deep learning models. Key observations from the preliminary results:

- **Class Imbalance Impact:** The model demonstrates better performance on the Negative class (F1-score: 0.84) compared to Non-Negative (F1-score: 0.71), reflecting the 1.72:1 class ratio present in the training data. This performance gap suggests the model has learned more robust patterns from the majority class.
- **Precision-Recall Trade-off:** The Negative class achieves high recall (0.86), indicating the model successfully identifies most negative emotions with few false negatives. However, Non-Negative recall is lower (0.68), suggesting approximately 32% of non-negative samples are misclassified as negative, likely due to class imbalance.
- **Computational Efficiency:** The SVM required only approximately 45 minutes for comprehensive hyperparameter optimization (grid search with 5-fold cross-validation) and final model training, demonstrating excellent computational efficiency suitable for resource-constrained environments.
- **Balanced Class Weights:** Employing balanced class weights (inversely proportional to class frequencies) improved performance on the minority class compared to

preliminary experiments with unweighted training, preventing the model from becoming biased toward the majority class.

C. Error Analysis

Preliminary analysis of misclassified samples reveals:

- 1) **Surprise Confusion:** Surprise samples (652 total) show the highest error rate, likely due to:
 - Underrepresentation in training data (5.4%)
 - Acoustic similarity to both fear (negative) and happy (non-negative)
- 2) **Neutral Ambiguity:** Some neutral samples are misclassified as negative, possibly due to:
 - Low emotional intensity making them difficult to distinguish
 - Variability in "neutral" expression across different actors
- 3) **Cross-Dataset Variability:** Differences in recording conditions and actor styles across the four source databases may contribute to misclassification.

D. Comparison with Literature

Our SVM baseline performance (79.57% accuracy) compares favorably with related work on emotion recognition:

- Nwe et al. [?] reported approximately 78% accuracy on 6-class emotion recognition using similar MFCC-based features and SVM classification.
- Zhang et al. [?] achieved comparable performance in the 75-80% range for binary valence (positive/negative) classification tasks.

These comparisons provide confidence that our baseline implementation is competitive with existing literature and that our dual feature extraction strategy (statistical aggregation for SVM, sequential representation for deep learning) is appropriate and effective for this binary emotion classification task.

V. REMAINING WORK AND UPDATED TIMELINE

A. Completed Tasks (Weeks 1-6)

Literature review and methodology design
 Dataset collection, merging, and preprocessing
 Dual feature extraction pipeline implementation
 SVM baseline development and evaluation
 CNN architecture design and initial implementation
 Progress report preparation

B. In-Progress Tasks (Week 6-7)

- CNN model training and hyperparameter tuning (80% complete)
- Initial performance evaluation and comparison with SVM
- Training dynamics visualization (loss curves, accuracy)

TABLE III
 REMAINING WORK SCHEDULE

Week	Task	Description
7	CNN Completion	Finalize training, evaluate, and analyze results
7-8	Bi-LSTM Implementation	Architecture implementation, training, and evaluation
8-9	Comprehensive Analysis	<ul style="list-style-type: none"> • Performance comparison • Statistical significance testing • Confusion matrix analysis • Training time comparison • Error pattern identification
9-10	Documentation	<ul style="list-style-type: none"> • Final report writing • Conference paper preparation • Code documentation • README and reproduction guide

C. Remaining Tasks (Week 7-10)

D. Expected Outcomes

Upon project completion, we expect to deliver:

- 1) **Final Report:** IEEE conference format with complete results
- 2) **Trained Models:** All three models with checkpoints
- 3) **Source Code:** Documented Python implementation with reproduction instructions
- 4) **Analysis:** Comprehensive comparison addressing:
 - Accuracy and per-class performance metrics
 - Training time and computational efficiency
 - Model complexity and deployment feasibility
 - Error patterns and failure cases
- 5) **Recommendations:** Practical guidelines for model selection based on application requirements

VI. CHALLENGES AND MITIGATION STRATEGIES

A. Encountered Challenges

- 1) **Class Imbalance:** The 1.72:1 ratio between Negative and Non-Negative classes presented a challenge. We addressed this by:
 - Using balanced class weights in SVM
 - Planning to use weighted loss functions for deep learning models
 - Ensuring stratified sampling in train-test split

- 2) **Surprise Underrepresentation:** With only 652 surprise samples (5.4%), the model has limited exposure to this emotion. Mitigation strategies include:
 - Careful analysis of surprise-specific errors
 - Consideration of data augmentation techniques (pitch shifting, time stretching)
 - Acknowledging this limitation in final analysis

- 3) **Cross-Dataset Variability:** Merging four databases with different recording conditions, actor demographics, and expression styles introduced variability. We addressed this through:

- Consistent preprocessing pipeline (resampling to 22,050 Hz)
- Normalization of extracted features
- Random shuffling before train-test split to mix sources

B. Anticipated Challenges

1) *Deep Learning Training Time*: Deep learning models (especially Bi-LSTM) may require extended training time on our hardware. We plan to:

- Use early stopping to prevent unnecessary epochs
- Implement efficient data loading with TensorFlow pipelines
- Consider reducing batch size if memory constraints arise

2) *Hyperparameter Optimization*: Deep learning models have many hyperparameters. We will:

- Start with literature-recommended values
- Perform limited grid search on critical parameters (learning rate, dropout)
- Use validation set performance for model selection

VII. CONCLUSION

This progress report documents significant advancement in our binary speech emotion recognition research. We have successfully:

- 1) Merged four publicly available databases into a comprehensive dataset of 12,162 audio samples
- 2) Implemented dual feature extraction strategies for traditional and deep learning approaches
- 3) Developed and evaluated an SVM baseline achieving 79.57% test accuracy
- 4) Initiated CNN implementation with architecture design completed

The SVM baseline demonstrates that binary emotion classification is feasible with traditional machine learning approaches, achieving balanced performance ($F1 = 0.80$) despite class imbalance. Preliminary error analysis has identified specific challenges, including surprise underrepresentation and neutral ambiguity, which will inform the development and evaluation of deep learning models.

With CNN training nearing completion and Bi-LSTM implementation scheduled for the next phase, we remain on track to complete all planned experiments and deliver comprehensive analysis by Week 10. The remaining work will focus on completing deep learning implementations, conducting thorough comparative analysis, and documenting findings in a final conference paper.

ACKNOWLEDGMENT

We thank the Department of Informatics Engineering, Universitas Syiah Kuala, for providing computational resources and guidance. We acknowledge the creators of CREMA-D, RAVDESS, SAVEE, and TESS for making their datasets publicly available. We also thank our supervisors for valuable feedback on this progress report.

REFERENCES

- 1 Schuller, B. W. and Batliner, A., "Speech emotion recognition: Two decades in a nutshell, benchmarks, and ongoing trends," *Communications of the ACM*, vol. 61, no. 5, pp. 90–99, 2018.
- 2 El Ayadi, M., Kamel, M. S., and Karray, F., "Survey on speech emotion recognition: Features, classification schemes, and databases," *Pattern Recognition*, vol. 44, no. 3, pp. 572–587, 2011.
- 3 Nwe, T. L., Foo, S. W., and De Silva, L. C., "Speech emotion recognition using hidden markov models," *Speech communication*, vol. 41, no. 4, pp. 603–623, 2003.
- 4 Badshah, A. M., Ahmad, J., Rahim, N., and Baik, S. W., "Speech emotion recognition from spectrograms with deep convolutional neural network," in *2017 International Conference on Platform Technology and Service (PlatCon)*. IEEE, 2017, pp. 1–5.
- 5 Sarma, M., Ghahremani, P., Povey, D., Goel, N. K., Sarma, K. K., and Dehak, N., "Emotion identification from raw speech signals using DNNs," in *Proc. Interspeech 2018*, 2018, pp. 3097–3101.
- 6 Mirsamadi, S., Barsoum, E., and Zhang, C., "Automatic speech emotion recognition using recurrent neural networks with local attention," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 2227–2231.
- 7 Huang, Z., Dong, M., Mao, Q., and Zhan, Y., "Speech emotion recognition using deep neural network and extreme learning machine," *Frontiers in neuroRobotics*, vol. 8, p. 34, 2014.
- 8 Valstar, M., Gratch, J., Schuller, B., Ringeval, F., Lalanne, D., Torres Torres, M., Scherer, S., Stratou, G., Cowie, R., and Pantic, M., "AVEC 2016: Depression, mood, and emotion recognition workshop and challenge," in *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*. ACM, 2016, pp. 3–10.
- 9 Cummins, N., Scherer, S., Krajewski, J., Schnieder, S., Epps, J., and Quatieri, T. F., "A review of depression and suicide risk assessment using speech analysis," *Speech Communication*, vol. 71, pp. 10–49, 2015.
- 10 Davis, S. and Mermelstein, P., "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," *IEEE transactions on acoustics, speech, and signal processing*, vol. 28, no. 4, pp. 357–366, 1980.
- 11 Eyben, F., Wöllmer, M., and Schuller, B., "openSMILE: the Munich versatile and fast open-source audio feature extractor," in *Proceedings of the 18th ACM International Conference on Multimedia*. ACM, 2010, pp. 1459–1462.
- 12 Cao, H., Cooper, D. G., Keutmann, M. K., Gur, R. C., Nenkova, A., and Verma, R., "CREMA-D: Crowd-sourced Emotional Multimodal Actors Dataset," *IEEE Transactions on Affective Computing*, vol. 5, no. 4, pp. 377–390, 2014.
- 13 Livingstone, S. R. and Russo, F. A., "The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English," *PloS one*, vol. 13, no. 5, p. e0196391, 2018.
- 14 Haq, S. and Jackson, P. J., "Speaker-dependent audio-visual emotion recognition," in *Proceedings of the 4th International Conference on Auditory-Visual Speech Processing (AVSP)*. Tangalooma, Australia, 2008, pp. 53–58.
- 15 Dupuis, K. and Pichora-Fuller, M. K., "Toronto Emotional Speech Set (TESS)," Scholars Portal Dataverse, 2010, available at: <https://tspace.library.utoronto.ca/handle/1807/24487>.
- 16 Zhang, S., Zhang, S., Huang, T., and Gao, W., "Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching," *IEEE Transactions on Multimedia*, vol. 20, no. 6, pp. 1576–1590, 2018.