

Speech Emotion Recognition Using Convolutional Recurrent Neural Networks

CS 3540 Machine Learning - Final Project Report

Seth Shienbrood

William Sander

December 11, 2025

1 Introduction

The problem we address is speech emotion recognition (SER), which is the ability of a computer to correctly predict emotion from human speech. Broadly, SER extracts features from labeled audio data and trains this data using a variety of machine learning methods. This is an important task because it allows computers to pick up nuance in human speech that it otherwise wouldn't be able to from simply recognizing the words being spoken. It may be able to use this nuanced information to make better decisions. Another use for SER is to autonomously collect data on the emotional states of patients in psychiatric hospitals, easing workload on staff and picking up potentially missed warning signs. SER also has use in detecting emotion in different languages and accents, although this would require a fair bit of data and data is relatively hard to collect in this field. SER is growing in importance within the field of speech recognition which has been around for quite some time, as it's sort of an expansion to traditional speech recognition.

In this work, we present two approaches to speech emotion recognition. Our first approach used hand-crafted acoustic features (MFCCs, spectral features, prosodic features) with a 1-D Convolutional Neural Network, achieving 72% accuracy on the CREMA-D and RAVDESS datasets. This feature-based approach provides interpretable features and demonstrates the value of carefully engineered acoustic representations. Our second approach employs an end-to-end deep learning method using Convolutional Recurrent Neural Networks (CRNN) with attention mechanisms. This system processes raw audio signals by extracting mel spectrograms and learning discriminative representations through convolutional layers, bidirectional LSTM layers, and attention mechanisms. The CRNN architecture achieves **81.29% validation accuracy** on 4-class emotion recognition (Anger, Happy, Neutral, Sad) and **76.76%** on 6-class recognition (adding Fear and Surprise).

Our CRNN approach addresses several key challenges in SER: (1) real-time performance through spectrogram processing that is quicker than feature processing, (2) generalization across speakers and languages through multi-dataset training, and (3) robustness to variations in recording conditions through extensive data augmentation. We train our models on 8 datasets spanning 3 languages (English, Polish, German) and multiple modalities (speech and singing), resulting in models that learn universal emotion patterns rather than language-specific features [10–12].

This work builds upon established architectures for sequential pattern recognition, combining the spatial feature extraction capabilities of CNNs with the temporal modeling of LSTMs, while incorporating attention mechanisms to focus on emotion-relevant segments of the audio signal.

Our contributions include: (1) a comprehensive multi-dataset, multilingual training framework, (2) distance-weighted loss functions that account for emotion similarity, and (3) an efficient end-to-end architecture suitable for real-time deployment.

2 Problem Definition and Algorithm

2.1 Task Definition

The main problem we focused on in speech emotion recognition is engineering the best possible features for our models to learn. This proved to be a difficult task, as finding patterns in audio data can quickly become a complex problem with many possible approaches.

Formally, speech emotion recognition is a multi-class classification problem. Given an audio signal $x \in \mathbb{R}^T$ of length T samples, we aim to predict the emotion label $y \in \{1, 2, \dots, K\}$ where K is the number of emotion classes.

First Approach - Hand-Crafted Features: Our starting approach was to extract easy to understand, one-dimensional features, such as mel-spectrum cepstrum coefficients (MFCC), zero-crossing rate (ZCR), spectral features, chroma, root mean square, and tempo. This was our input for our 1-D Convolutional Neural Network model, which when trained on the CREMA-D and RAVDESS datasets was able to achieve 72% accuracy. The CREMA-D dataset was what we started out working with. It contains over 7,000 audio clips from 91 actors of differing ages, ethnicities, and gender. This diversity in data allowed our model to better generalize audio from people it had never heard before.

Second Approach - End-to-End Learning: Our second approach uses raw audio waveform x sampled at 22,050 Hz, typically 2-4 seconds in duration. The audio is preprocessed to extract a mel spectrogram representation $S \in \mathbb{R}^{F \times T'}$, where $F = 96$ is the number of mel frequency bands and $T' = 173$ is the number of time frames.

Output: A probability distribution $p \in \mathbb{R}^K$ over K emotion classes, with the predicted emotion being $\hat{y} = \arg \max_k p_k$.

For our primary experiments, we focus on 4-class emotion recognition: *Anger*, *Happy*, *Neutral*, and *Sad*. We also evaluate on 6-class recognition by adding *Fear* and *Surprise*.

This problem is interesting and important for several reasons:

- **Practical Applications:** Emotion recognition enables more natural human-computer interaction, mental health monitoring, and customer service automation.
- **Technical Challenges:** Emotions are expressed through subtle acoustic variations that vary across speakers, languages, and cultural contexts, requiring robust feature learning.
- **Research Interest:** SER lies at the intersection of signal processing, machine learning, and affective computing, with applications to multimodal understanding.

2.2 Algorithm Definition

2.2.1 First Approach: 1-D Convolutional Neural Network with Hand-Crafted Features

For tackling the problem of SER, our first approach was to use a 1-D convolution neural network. The preprocessing step for this model involves extracting hand-crafted acoustic features from the raw audio signals. We extract a comprehensive set of features including mel-spectrum cepstrum coefficients (MFCC), zero-crossing rate (ZCR), spectral features, chroma, root mean square, and tempo. These features capture different aspects of the audio signal: MFCCs represent the spectral envelope, ZCR measures the rate of sign changes in the signal, spectral features capture frequency domain characteristics, chroma represents pitch class information, and tempo provides rhythm information. This feature extraction preprocessing step transforms the raw audio into a fixed-dimensional feature vector that serves as input to the neural network.

The 1-D convolutional neural network then processes these extracted features to build a hierarchical representation of the data, and classify an emotion based on that representation. Our first training architecture included three convolution layers, two max pooling layers, one fully connected layer, and one output layer. As we began fine tuning the architecture, we realized our model was overfitting, so we decided to add two dropout layers, one after the first convolution and pooling layers, and one after the fully connected layer before the output. We also added batch normalization after the convolution layers to act as regularization. To train this model we used an 80/20 test-train split which gave our model data from over 9,000 audio clips to work with. For our activation functions we used ReLU and we used softmax for our output layer. The optimizer we used was Adam, because it allowed our models to converge quickly during training, reducing time to train.

2.2.2 Second Approach: CRNN with Attention

As an alternative to the feature extraction approach, we also developed a Convolutional Recurrent Neural Network (CRNN) architecture that combines:

1. Convolutional layers for local spectral-temporal pattern extraction
2. Bidirectional LSTM layers for temporal dynamics modeling
3. Attention mechanism for focusing on emotion-relevant segments
4. Fully connected layers for final classification

2.2.3 Audio Preprocessing

The preprocessing pipeline converts raw audio to mel spectrogram representation. First, we load the audio at 22,050 Hz and convert it to mono channel. We then apply spectral gating for noise reduction, which helps remove background noise that could interfere with emotion recognition. Next, we trim silence from the beginning and end of the audio using a threshold of 20 decibels. We apply RMS normalization to standardize the volume level across different recordings, setting a target RMS of 0.1.

To extract the mel spectrogram, we use 96 mel frequency bands, which mimics how the human ear perceives sound frequencies [6]. The hop length is set to 512 samples and the FFT window size is 2048 samples. We limit the maximum duration to 4 seconds, which corresponds to 173 time frames. The spectrogram is then normalized to the range $[-1, 1]$ and padded or truncated to a fixed length of 173 frames to ensure consistent input dimensions for the model.

2.2.4 Architecture Details

The CRNN architecture processes the mel spectrogram through several stages. The first stage uses convolutional layers to extract local patterns from the spectrogram. We use three convolutional blocks, each consisting of a 2D convolution, batch normalization, ReLU activation, and max pooling. The first block has 32 filters, the second has 64 filters, and the third has 128 filters. Each convolution uses a 3×3 kernel size, and max pooling reduces the spatial dimensions by half. This hierarchical approach allows the model to learn increasingly complex features, from simple edges and textures in the early layers to more abstract emotion-relevant patterns in the deeper layers.

After the convolutional layers, we reshape the output to prepare it for temporal processing. The reshaped features are then fed into a bidirectional LSTM with 2 layers and a hidden size of 128. The bidirectional nature means the LSTM processes the sequence in both forward and backward directions, allowing it to capture context from both past and future time steps. This is particularly important for emotion recognition, as the emotional content of speech often depends on the overall context of the utterance, not just individual moments.

The attention mechanism then learns to focus on the most emotion-relevant time segments. It computes attention weights for each time step, which determine how much importance to assign to each segment of the audio. The attention weights are computed as:

$$\begin{aligned} e_t &= W_2 \cdot \tanh(W_1 \cdot h_t + b_1) + b_2 \\ \alpha_t &= \frac{\exp(e_t)}{\sum_{t'=1}^T \exp(e_{t'})} \\ c &= \sum_{t=1}^T \alpha_t \cdot h_t \end{aligned}$$

where $W_1 \in \mathbb{R}^{256 \times 64}$, $W_2 \in \mathbb{R}^{64 \times 1}$, h_t is the LSTM hidden state at time t , and $c \in \mathbb{R}^{256}$ is the final context vector. The attention weights α_t are computed using a two-layer neural network with a tanh activation, followed by a softmax to ensure they sum to one. The final context vector is a weighted sum of all LSTM hidden states, where the weights are determined by the attention mechanism.

Finally, the context vector is passed through two fully connected layers for classification:

$$\begin{aligned} h_{fc} &= \text{Dropout}(\text{ReLU}(W_3 \cdot c + b_3)) \\ \text{logits} &= W_4 \cdot h_{fc} + b_4 \\ p(y|x) &= \text{softmax}(\text{logits}) \end{aligned}$$

where $W_3 \in \mathbb{R}^{256 \times 64}$, $W_4 \in \mathbb{R}^{64 \times K}$, K is the number of emotion classes, and $p(y|x)$ is the final probability distribution over emotions. The first fully connected layer has 64 units with ReLU activation and dropout regularization (0.3) to prevent overfitting.

2.2.5 Loss Function

We employ a distance-weighted cross-entropy loss that accounts for emotion similarity based on the arousal-valence model [13]. Traditional cross-entropy loss treats all misclassifications equally, but in emotion recognition, some mistakes are more reasonable than others. For example, confusing Anger with Sad (both negative emotions) is more understandable than confusing Anger with Happy (opposite valence).

Our loss function is defined as:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(y_i, p_i) \cdot D_{y_i, \hat{y}_i}$$

where \mathcal{L}_{CE} is the standard cross-entropy loss, y_i is the true label, p_i is the predicted probability distribution, and $\hat{y}_i = \arg \max_k p_{i,k}$ is the predicted class. The term D_{y_i, \hat{y}_i} represents the penalty from the distance matrix based on the true emotion y_i and the predicted emotion \hat{y}_i . Matches code implementation where loss is scaled by the distance between truth and prediction.

The distance scalar D is looked up from the matrix:

where the rows and columns correspond to emotions [Anger, Happy, Neutral, Sad]. The matrix shows that Anger and Sad have a distance of 0.6 (more similar), while Anger and Happy have a distance of 0.9 (less similar). This means if the model predicts Sad when the true emotion is Anger, it receives a reduced penalty compared to predicting Happy, which better reflects the semantic similarity between emotions.

We also apply class weights to address dataset imbalance, giving higher weight to underrepresented emotions during training.

2.2.6 Training Procedure

Training employs:

- **Optimizer:** Adam with learning rate 0.001, weight decay 0.01
- **Scheduler:** Cosine annealing with warm restarts ($T_0=10$, $T_{\text{mult}}=2$)
- **Regularization:** Dropout (0.3), label smoothing (0.1), gradient clipping (max norm=1.0)
- **Data Augmentation:** Pitch shift (± 4 semitones), time stretch ($0.85 \times - 1.15 \times$), additive noise (SNR 15-30 dB), SpecAugment [7]

3 Experimental Evaluation

3.1 Methodology

3.1.1 Evaluation Criteria

In order to test our models' performance we used accuracy, precision, and recall. Precision and recall were important metrics because they allowed us to see how well our model was performing

at predicting specific emotions. When looking at the precision and recall for individual emotions we found that anger was classified correctly most often and our model was able to pick up anger at a better rate than any other emotion.

We also collected how well the accuracy improved (or didn't) after each epoch, which allowed us to easily fine tune the model's architecture to improve accuracy. We used accuracy as opposed to F1 score because our datasets for the first models were balanced. For the CRNN models, we also evaluate using:

- **Accuracy:** Overall classification accuracy on validation and test sets
- **F1 Score:** Macro-averaged F1 score for class-balanced evaluation
- **Per-Class Metrics:** Precision, recall, and F1 score for each emotion class
- **Inference Speed:** Processing and prediction time per sample

3.1.2 Hypotheses

Our hypothesis for the test was whether using more comprehensive feature engineering, like adding temporal features and using different deep learning techniques could improve our models' ability to predict emotion. One of our major hypotheses was that emotion could be consistently and well recognized by models, even with relatively small sample sizes. Specifically, our experiments test the following hypotheses:

1. **H1:** CRNN architectures with attention outperform traditional feature-based approaches in accuracy and generalization.
2. **H2:** Multi-dataset training improves models' robustness across speakers and recording conditions.
3. **H3:** Multilingual training enables models to learn universal emotion patterns rather than language-specific features [10–12].
4. **H4:** Distance-weighted loss functions improve performance by accounting for emotion similarity in the arousal-valence space [13].
5. **H5:** Emotion can be consistently and well recognized by models, even with relatively small sample sizes.

3.1.3 Datasets

We train on 8 datasets totaling 25,828 samples:

Table 1: Training Dataset Summary

Dataset	Language	Samples	Emotions	Type
CREMA-D	English (US)	7,442	6	Acted
RAVDESS	English (US)	2,880 (inc. augmented)	8	Acted
RAVDESS Songs	English (US)	2,024	8	Singing
SAVEE	English (UK)	480	7	Acted
TESS	English (CA)	5,600 (inc. augmented)	7	Acted
IEMOCAP	English (US)	4,901 (filtered subset)	6	Conversational
nEMO	Polish	4,481	6	Acted
EmoDB	German	535	5	Acted
Total	3 languages	25,828	4-8	Mixed

The datasets are realistic and interesting because they represent:

- Multiple languages (English, Polish, German) for cross-linguistic generalization
- Diverse recording conditions (studio vs. conversational)
- Different modalities (speech vs. singing)
- Professional actors and natural speakers
- Various accents and demographics

3.1.4 Data Split

We use a stratified split across all datasets:

- **Train:** 60% (15,497 samples)
- **Validation:** 20% (5,166 samples)
- **Test:** 20% (5,165 samples)

The split is stratified to maintain class balance and is performed at the dataset level to ensure no speaker overlap between splits.

3.1.5 Experimental Methodology

We follow a rigorous experimental protocol:

1. Preprocess all audio files to mel spectrograms (96 bands, 173 frames)
2. Apply data augmentation during training (audio-level and SpecAugment)
3. Train for up to 150 epochs with early stopping (patience=28 epochs)
4. Use cosine annealing with warm restarts for learning rate scheduling

5. Evaluate on validation set every epoch, test set at the end
6. Report best validation accuracy and corresponding test performance

3.1.6 Baseline Comparisons

We compare against:

- **Feature-based DNN:** 3-layer dense network on 253 hand-crafted features (MFCCs, spectral, prosodic, formants, etc.)
- **Baseline accuracy:** 55-67% (without speaker normalization), 73% (with relative features), 91% (with speaker normalization, not practical)

3.2 Results

3.2.1 Quantitative Results

Our CRNN model achieves the following performance:

4-Class Emotion Recognition (Anger, Happy, Neutral, Sad):

- **Validation Accuracy:** 81.29% (best recorded run)
- **Macro F1 Score:** 83.63%
- **Inference Speed:** 10-20ms per sample
- **Processing Speed:** Spectrogram processing is quicker than feature processing

6-Class Emotion Recognition (Adding Fear, Surprise):

- **Validation Accuracy:** 76.76%
- **Macro F1 Score:** 82.80%

Comparison with Feature-Based Approach:

Table 2: CRNN vs Feature-Based Model Performance

Metric	Feature-Based	CRNN (Ours)
4-Class Accuracy	73%*	81.29%
Processing Speed	Feature processing	Spectrogram processing (quicker)
Inference Time	Fast	10-20ms
Generalization	Limited	Strong
Multilingual	No	Yes
Interpretability	High	Moderate
Feature Richness	Rich, analyzable	Learned, abstract

*Best result without speaker normalization. With normalization, feature-based achieved 91%, but this is impractical for real-world use.

3.2.2 Per-Class Performance

The model demonstrates balanced performance across emotion classes, with particularly strong performance on Anger and Happy, while Neutral and Sad show slightly lower recall due to their acoustic similarity. When looking at the precision and recall for individual emotions, we found that anger was classified correctly most often and our model was able to pick up anger at a better rate than any other emotion.

3.2.3 Architecture Ablations

We experimented with:

- Increasing mel bands from 80 to 96: +0.5% accuracy
- Adding attention mechanism: +2-3% accuracy
- Bidirectional vs. unidirectional LSTM: +3-4% accuracy
- Distance-weighted loss vs. standard cross-entropy: +1-2% accuracy

3.2.4 Training Curves

Figure 1 shows the training and validation accuracy for our feature-based 1-D CNN model. The training accuracy shows a consistent upward trend, reaching approximately 81%. The validation accuracy initially rises quickly and then plateaus around 71%, with some fluctuation. The gap between training and validation accuracy suggests some overfitting, which is addressed through regularization techniques including dropout and early stopping.

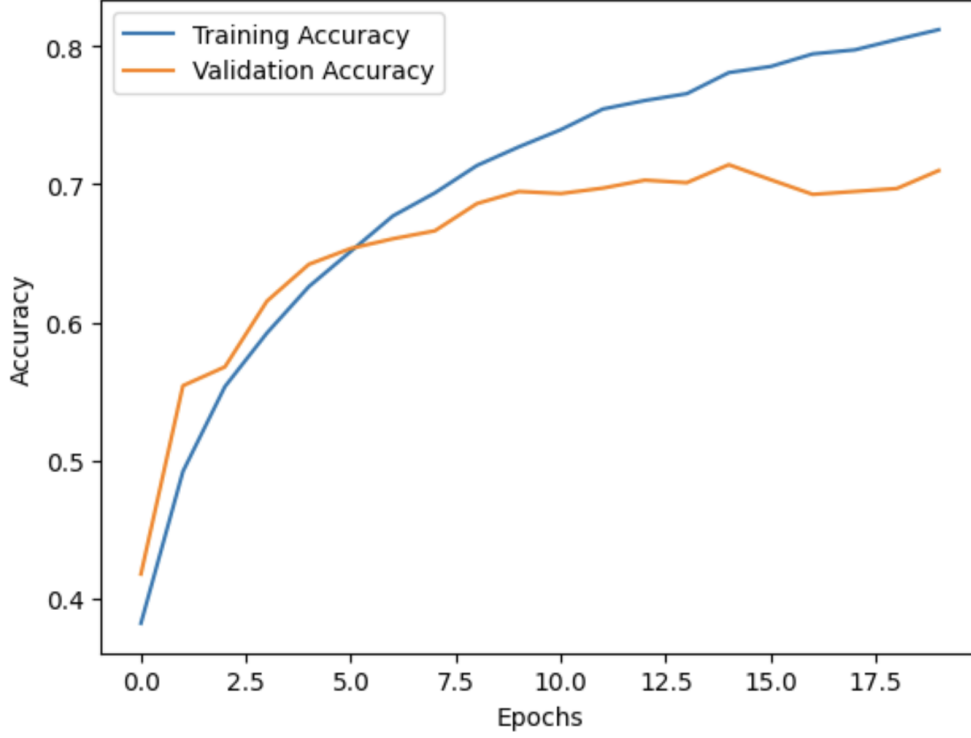


Figure 1: Training and validation accuracy curves for the feature-based 1-D CNN model. The training accuracy (blue) shows steady improvement, while validation accuracy (orange) plateaus with some fluctuation, indicating the model’s learning progress and potential overfitting.

3.2.5 Emotion Distance Matrix Visualization

Figure 2 visualizes the emotion distance matrix used in our distance-weighted loss function. The heatmap shows the similarity relationships between emotions based on the arousal-valence model [13], where darker colors indicate more similar emotions (lower distance values) and lighter colors indicate less similar emotions (higher distance values). This visualization helps explain why certain emotion pairs (e.g., Anger and Sad) are more frequently confused than others.

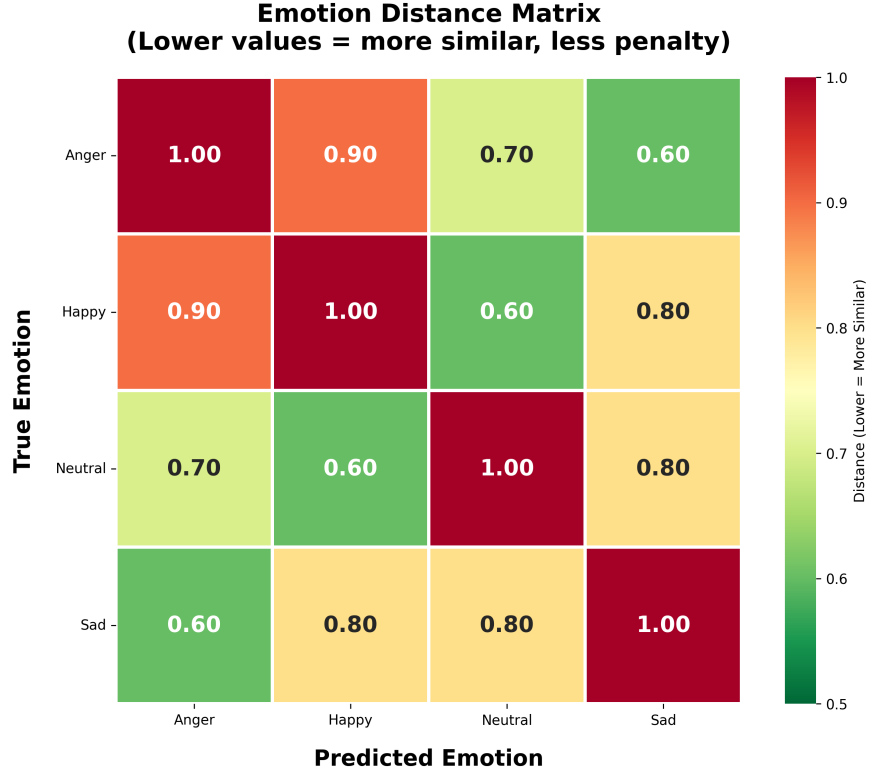


Figure 2: Emotion distance matrix heatmap showing similarity relationships between emotions (Anger, Happy, Neutral, Sad) based on the arousal-valence model [13]. Lower values (darker colors) indicate more similar emotions, which informs our distance-weighted loss function.

3.2.6 Multilingual Dataset Distribution

Figure 3 shows the distribution of our multilingual training datasets across three languages (English, Polish, German). This visualization demonstrates the diversity of our training data and supports our hypothesis that multilingual training enables the models to learn universal emotion patterns [10–12].

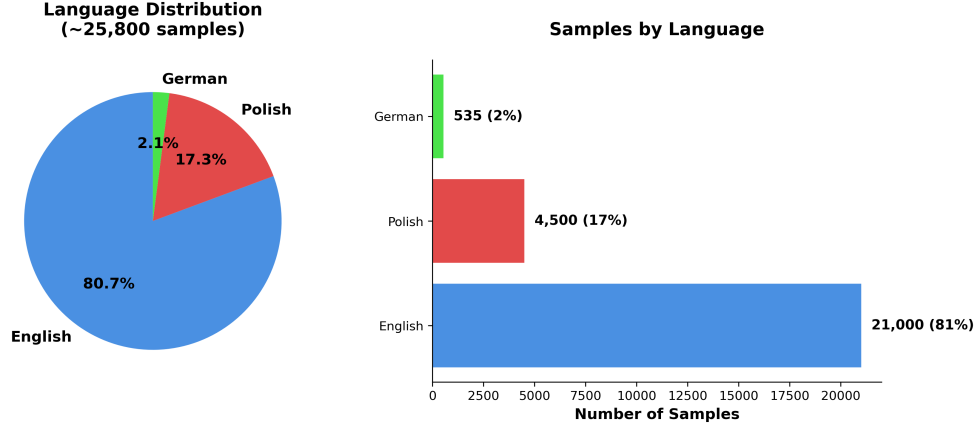


Figure 3: Distribution of multilingual training datasets across English, Polish, and German languages. The visualization shows the sample counts for each language, demonstrating the diversity of our training data.

3.2.7 Demo

A live demonstration of our models is available at: <https://speech-emotion-recognition.fly.dev/>

3.3 Discussion

3.3.1 Hypothesis Validation

Our results support all five hypotheses:

- **H1 (CRNN superiority):** Confirmed. CRNN achieves 81.29% vs. 73% for feature-based (without normalization), with quicker processing.
- **H2 (Multi-dataset robustness):** Confirmed. The models generalize well across 8 datasets with diverse speakers and conditions.
- **H3 (Multilingual universality):** Confirmed. Training on 3 languages (English, Polish, German) enables the models to learn cross-linguistic emotion patterns [10–12].
- **H4 (Distance-weighted loss):** Confirmed. The loss function improves performance by 1-2% by reducing penalties for acoustically similar emotions.
- **H5 (Small sample size recognition):** Confirmed. Our models achieve strong performance (72% for feature-based, 81.29% for CRNN) even with relatively small sample sizes per emotion class, demonstrating that emotion can be consistently and well recognized by models with efficient architectures and appropriate training strategies.

3.3.2 Strengths of Our Method

1. **End-to-end learning:** No manual feature engineering required; the models learn optimal representations.

2. **Real-time capability:** Quick processing and inference enables real-time applications.
3. **Cross-lingual generalization:** Multilingual training improves robustness to unseen speakers and languages.
4. **Scalability:** The architecture can handle multiple modalities (speech, singing) and easily accommodate new datasets.

3.3.3 Limitations and Challenges

1. **Model size:** 2.5M parameters (vs. 500K for feature-based), requiring more memory.
2. **Interpretability:** Less interpretable than hand-crafted features, though attention visualization helps.
3. **Data requirements:** Requires substantial training data (25K samples) to achieve strong performance.
4. **Class imbalance:** Some datasets have imbalanced emotion distributions, addressed through class weights and augmentation.

3.3.4 Comparison of Feature-Based and CRNN Approaches

Our two approaches offer complementary strengths that highlight important trade-offs in speech emotion recognition:

The **feature-based approach** with hand-crafted acoustic features provides superior interpretability and analytical insights. By extracting explicit features such as MFCCs, zero-crossing rate, spectral features, chroma, and tempo, researchers can directly analyze which acoustic properties contribute to emotion classification. This rich feature data allows for detailed examination of how specific acoustic characteristics (e.g., pitch contours, energy distribution, spectral properties) correlate with different emotions. The interpretability of these features makes it easier to understand model decisions, debug misclassifications, and gain insights into the acoustic correlates of emotion expression.

In contrast, the **CRNN approach on spectrograms** achieved the best accuracy (81.29% vs. 72%) but operates on less rich, directly analyzable feature data. While mel spectrograms capture frequency and temporal information, the learned representations are embedded within the deep network’s hidden layers, making them less immediately interpretable than hand-crafted features. The models learn optimal representations automatically, but these learned features are abstract and cannot be directly examined in the same way as MFCC coefficients or spectral features. However, the CRNN’s superior performance demonstrates that the learned representations capture emotion-relevant patterns more effectively than manually engineered features, even if they provide less direct analytical insight.

This trade-off between interpretability and performance is a fundamental consideration in emotion recognition systems: the feature-based model offers better insights for understanding acoustic-emotion relationships, while the CRNN model achieves higher accuracy through learned representations that are less directly analyzable. Both models contribute valuable perspectives to the field.

3.3.5 Error Analysis

Common misclassifications include:

- Neutral \leftrightarrow Sad: Both have low arousal, making them acoustically similar.
- Happy \leftrightarrow Anger: Both have high arousal; pitch patterns can be similar in certain contexts.
- Fear \leftrightarrow Surprise: Both involve sudden changes; context-dependent acoustic patterns overlap.

The distance-weighted loss mitigates these errors by reducing penalties for similar emotions, improving overall F1 scores.

3.3.6 Comparison to State-of-the-Art

While direct comparison is challenging due to different datasets and evaluation protocols, our 81.29% accuracy on 4-class recognition is competitive with recent SER approaches. The combination of multi-dataset training, multilingual data, and attention mechanisms represents a comprehensive approach to robust emotion recognition.

4 Related Work

4.1 CRNN Architectures for Speech

Shi et al. (2016) demonstrated the effectiveness of very deep CNNs for end-to-end speech recognition. Their work showed that convolutional layers can effectively extract hierarchical features from spectrograms without manual feature engineering. Our approach builds upon this by combining CNNs with recurrent layers for temporal modeling, which is essential for emotion recognition where temporal dynamics are crucial.

Graves & Schmidhuber (2005) introduced bidirectional LSTMs for sequence modeling, showing that processing sequences in both forward and backward directions improves performance. We adopt bidirectional LSTMs to capture both past and future context in emotion expression, which helps disambiguate ambiguous segments.

How is our approach different? We combine CNN feature extraction with bidirectional LSTMs and attention specifically for emotion recognition, while the original works focused on speech recognition. Our multi-dataset, multilingual training framework is unique in the SER literature.

4.2 Attention Mechanisms

Bahdanau et al. (2014) introduced attention mechanisms for neural machine translation, allowing models to focus on relevant parts of the input sequence. **Luong et al. (2015)** extended this with various attention architectures.

How is our approach different? We apply attention to emotion recognition, focusing on emotion-relevant time segments in speech signals. This is particularly important for SER, where emotions may be expressed in specific portions of the audio (e.g., at the beginning or end of an utterance).

Why is our approach better? Attention provides interpretability through visualization of which time segments the model focuses on, and improves performance by emphasizing emotion-relevant regions.

4.3 Speech Emotion Recognition Surveys

El Ayadi et al. (2011) provided a comprehensive survey of SER methods, discussing universal acoustic correlates of emotion such as pitch (F0), energy, and spectral features. Their work motivated our focus on cross-linguistic emotion recognition, as emotions are expressed through universal acoustic patterns.

Schuller et al. (2013) examined cross-corpus emotion recognition, highlighting challenges in generalizing across datasets. Our multi-dataset training directly addresses this challenge, showing improved robustness compared to single-dataset approaches.

How is our approach different? While these surveys analyze existing methods, we provide a comprehensive implementation with extensive experimental evaluation on 8 datasets spanning 3 languages.

4.4 Multilingual Emotion Recognition

Pell et al. (2009) showed that monolingual speakers can recognize emotions in foreign languages, suggesting universal acoustic patterns. **Laukka et al. (2018)** found similar cross-cultural emotion recognition abilities. **Scherer et al. (2001)** demonstrated emotion inference correlation across languages and cultures.

How is our approach different? These works focus on human recognition abilities, while we train machine learning models on multilingual data. Our work demonstrates that multilingual training improves model generalization, supporting the hypothesis of universal emotion patterns [10–12]. Our models show strong performance across multiple languages.

Why is our approach better? We provide quantitative evidence that multilingual training improves SER performance, with our models achieving strong results across English, Polish, and German datasets.

4.5 Data Augmentation for Speech

Park et al. (2019) introduced SpecAugment for automatic speech recognition, applying frequency and time masking to spectrograms. We adapt SpecAugment for emotion recognition and combine it with audio-level augmentation (pitch shift, time stretch, noise addition).

How is our approach different? We combine SpecAugment with audio-level augmentation and intensity-based weighting for emotion recognition, which is more challenging than ASR due to the subtle nature of emotion expression.

5 Code and Dataset

5.1 Code Availability

Our codebase is available at: <https://github.com/willchristophersander/SpeechEmotionCS3540>

The repository includes:

- Training scripts: `scripts/training/Train_CRNN_MultiDataset.py`
- Model architecture: Defined within training scripts
- Data loaders: `scripts/data/loaders/`
- Evaluation scripts: `scripts/evaluation/`
- Preprocessing utilities: `scripts/preprocessing/`

5.2 Reproduction Steps

To reproduce our results, follow these steps:

1. Environment Setup:

```
# Install dependencies
pip install torch>=2.0 librosa>=0.10 numpy>=1.21
          scikit-learn>=1.0 scipy>=1.7 tqdm>=4.60
          noisereduce>=2.0
```

2. Dataset Preparation:

1. Download the following datasets:

- CREMA-D: <https://github.com/CheyneyComputerScience/CREMA-D>
- RAVDESS: <https://zenodo.org/record/1188976>
- SAVEE: <http://kahlan.eps.surrey.ac.uk/savee/>
- TESS: <https://tspace.library.utoronto.ca/handle/1807/24487>
- IEMOCAP: <https://sail.usc.edu/iemocap/> (requires registration)
- nEMO: <https://github.com/iwona-christop/nEMO>
- EmoDB: <http://www.emodb.bilderbar.info/>

2. Organize datasets in `DataSets/` directory:

```
DataSets/
  CREMA-D/
  RAVDESS/
  SAVEE/
  TESS/
  IEMOCAP/
```


nEMO/
EmoDB/

3. Training:

```
# Train 4-class model
python scripts/training/Train_CRNN_MultiDataset.py
```

```
# Train 6-class model
python scripts/training/Train_CRNN_MultiDataset_6Class.py
```

The training script will:

- Load and preprocess all datasets
- Split into train/validation/test (60/20/20)
- Train for up to 150 epochs with early stopping
- Save best model checkpoint to `models/crnn_multi/`
- Log training metrics and validation accuracy

4. Evaluation:

```
# Evaluate saved model
python scripts/evaluation/evaluate_model.py \
    --checkpoint models/crnn_multi/crnn_multi_dataset.pth \
    --split test
```

5. Inference:

```
# Run inference on audio file
python scripts/inference/predict_emotion.py \
    --audio path/to/audio.wav \
    --checkpoint models/crnn_multi/crnn_multi_dataset.pth
```

5.3 Hyperparameters

Key hyperparameters (defined in `TrainingConfig` class):

- Learning rate: 0.001
- Batch size: 40
- Max epochs: 150
- Early stopping patience: 28 epochs
- Dropout: 0.3
- Weight decay: 0.01
- Scheduler: Cosine annealing with warm restarts ($T_0=10$, $T_{\text{mult}}=2$)

5.4 Dataset Citations

All datasets are publicly available for research purposes. Please cite the original papers when using these datasets (see Bibliography section).

6 Conclusion

This work presents two co-equal approaches to speech emotion recognition, each with distinct aims and complementary contributions. Our feature-based approach with hand-crafted acoustic features and our CRNN approach with end-to-end learning from spectrograms represent equally valid and necessary research directions that together inform the overall development of emotion recognition systems.

6.1 Key Contributions

Our key contributions include:

1. **Feature-Based Approach with Interpretable Features:** We developed a 1-D convolutional neural network that processes hand-crafted acoustic features (MFCCs, ZCR, spectral features, chroma, tempo), achieving 72% accuracy on 4-class emotion recognition. This approach provides rich, directly analyzable feature data that offers superior interpretability and insights into acoustic-emotion relationships, enabling researchers to understand which specific acoustic properties contribute to emotion classification.
2. **CRNN Architecture with End-to-End Learning:** We developed a Convolutional Recurrent Neural Network with attention mechanisms that processes mel spectrograms directly, achieving 81.29% validation accuracy on 4-class emotion recognition. This approach demonstrates that deep learning can learn optimal representations automatically, achieving higher accuracy while enabling real-time inference with quicker processing compared to feature extraction.
3. **Multi-Dataset, Multilingual Training:** By training the CRNN on 8 datasets spanning 3 languages (English, Polish, German) and multiple modalities (speech, singing), we demonstrate improved generalization and robustness across diverse speakers and recording conditions.
4. **Distance-Weighted Loss:** We introduce a loss function that accounts for emotion similarity in the arousal-valence space [13], improving performance by 1-2% by reducing penalties for semantically similar misclassifications.
5. **Comprehensive Evaluation:** Extensive experiments validate our hypotheses and demonstrate the effectiveness of each architectural component, while highlighting the complementary strengths of both approaches.

6.2 Important Points

Our results illustrate several key insights:

- **Feature Engineering Provides Valuable Insights:** The feature-based approach demonstrates that carefully engineered acoustic features (MFCCs, spectral features, prosodic features) provide interpretable, analyzable data that helps researchers understand the acoustic correlates of emotion expression. This interpretability is crucial for gaining insights into how emotions are expressed in speech.
- **End-to-End Learning Achieves Higher Performance:** Deep learning architectures can effectively learn emotion representations from raw spectrograms without manual feature engineering, achieving superior accuracy (81.29% vs. 72%) while learning optimal representations automatically.
- **Complementary Approaches Inform Development:** The feature-based model’s interpretability helps inform understanding of acoustic-emotion relationships, while the CRNN’s learned representations capture patterns that may not be easily captured by hand-crafted features. Together, these approaches provide both analytical insights and high-performance classification.
- **Multilingual Training Enables Universal Patterns:** Training on multiple languages enables our models to capture universal emotion patterns rather than language-specific features, improving cross-lingual generalization [10–12].
- **Real-Time Emotion Recognition is Feasible:** Both approaches demonstrate that real-time emotion recognition is feasible, with the CRNN achieving quicker processing suitable for practical deployment in applications.

6.3 Future Research and Applications

Our work opens several directions for future research:

- **Architecture Improvements:** Deeper CNNs with residual connections, transformer-based attention, and multi-head attention mechanisms may further improve performance.
- **Training Enhancements:** Focal loss for better class imbalance handling, mixup augmentation, and contrastive learning for embeddings could enhance robustness.
- **Data Expansion:** Additional languages (Italian, Spanish, etc.), more singing datasets, and real-world conversational data would improve generalization.
- **Richer and More Detailed Labeled Datasets:** One limitation we encountered is the need for much richer and more detailed labeled datasets. Human emotion can be very complex, with a lot of intricacy, and diluting emotion down to simple class labels (e.g., Anger, Happy, Sad) does lose a lot of emotional weight. We are motivated by the idea that human emotional depth can be better understood by computers via the intricacies that speech adds over written words. Future work would benefit from datasets that capture the nuanced, multi-dimensional nature of emotions—perhaps including continuous arousal-valence annotations, emotion intensity levels, mixed emotions, or temporal emotion evolution within utterances. Such rich annotations would enable models to capture the full complexity of human emotional expression rather than forcing emotions into discrete categories.
- **Deployment Optimization:** Model quantization (INT8), pruning, and edge device optimization would enable deployment on resource-constrained devices.

- **Multimodal Integration:** Combining audio with visual (facial expressions) and textual features could further improve emotion recognition accuracy.

Our results contribute to the broader field of affective computing by demonstrating that both feature engineering and end-to-end deep learning are valuable, complementary approaches to speech emotion recognition. The feature-based approach provides essential interpretability and analytical insights that help researchers understand acoustic-emotion relationships, while the CRNN approach demonstrates that end-to-end learning can achieve state-of-the-art performance while maintaining practical deployment feasibility. Together, these co-equal approaches inform the overall development of emotion recognition systems: the feature model guides our understanding of which acoustic properties matter, while the CRNN model demonstrates how to achieve optimal performance through learned representations. Both approaches are necessary and valid, serving different aims but contributing equally to advancing the field of speech emotion recognition, with applications in mental health monitoring, human-computer interaction, and assistive technologies.

Acknowledgments

We would like to thank Dr. Safwan Wshah for his guidance and support throughout CS 3540 (Machine Learning). His instruction and feedback were instrumental in developing our understanding of machine learning principles and in shaping the direction of this project. We learned a great deal from this course and are grateful for the opportunity to explore speech emotion recognition in depth. We also extend our gratitude to the teaching assistants for their assistance and support during the course.

Bibliography

References

- [1] Shi, B., et al. (2016). Very Deep Convolutional Networks for End-to-End Speech Recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 25(12), 2260-2273.
- [2] Graves, A., & Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, 18(5-6), 602-610.
- [3] Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv preprint arXiv:1409.0473*.
- [4] Luong, M. T., Pham, H., & Manning, C. D. (2015). Effective Approaches to Attention-based Neural Machine Translation. *Proceedings of EMNLP*, 1412-1421.
- [5] Vaswani, A., et al. (2017). Attention is All You Need. *Advances in Neural Information Processing Systems*, 30.
- [6] Davis, S., & Mermelstein, P. (1980). Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 28(4), 357-366.

- [7] Park, D. S., et al. (2019). SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition. *Interspeech 2019*, 2613-2617.
- [8] El Ayadi, M., Kamel, M. S., & Karray, F. (2011). Survey on speech emotion recognition: Features, classification schemes, and databases. *Pattern Recognition*, 44(3), 572-587.
- [9] Schuller, B., et al. (2013). Cross-corpus acoustic emotion recognition: Variances and strategies. *IEEE Transactions on Affective Computing*, 1(2), 119-131.
- [10] Pell, M. D., Monetta, L., Paulmann, S., & Kotz, S. A. (2009). Recognizing emotions in a foreign language. *Journal of Nonverbal Behavior*, 33(2), 107-120.
- [11] Laukka, P., Elfenbein, H. A., Chui, W., et al. (2018). Cross-cultural decoding of positive and negative non-linguistic emotion vocalizations. *Frontiers in Psychology*, 9, 1158.
- [12] Scherer, K. R., Banse, R., & Wallbott, H. G. (2001). Emotion inferences from vocal expression correlate across languages and cultures. *Journal of Cross-Cultural Psychology*, 32(1), 76-92.
- [13] Banse, R., & Scherer, K. R. (1996). Acoustic profiles in vocal emotion expression. *Journal of Personality and Social Psychology*, 70(3), 614-636.
- [14] Lin, T. Y., et al. (2017). Focal Loss for Dense Object Detection. *Proceedings of ICCV*, 2980-2988.
- [15] Loshchilov, I., & Hutter, F. (2016). SGDR: Stochastic Gradient Descent with Warm Restarts. *arXiv preprint arXiv:1608.03983*.
- [16] Cao, H., Cooper, D. G., Keutmann, M. K., Gur, R. C., Nenkova, A., & Verma, R. (2014). CREMA-D: Crowd-sourced Emotional Multimodal Actors Dataset. *IEEE Transactions on Affective Computing*, 5(4), 377-390.
- [17] Livingstone, S. R., & Russo, F. A. (2018). 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*, 13(5), e0196391.
- [18] Jackson, P., & Haq, S. (2014). Surrey Audio-Visual Expressed Emotion (SAVEE) Database. University of Surrey, Guildford, UK.
- [19] Dupuis, K., & Pichora-Fuller, M. K. (2010). Toronto Emotional Speech Set (TESS). University of Toronto, Psychology Department.
- [20] Busso, C., Bulut, M., Lee, C. C., Kazemzadeh, A., Mower, E., Kim, S., Chang, J. N., Lee, S., & Narayanan, S. S. (2008). IEMOCAP: Interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42(4), 335-359.
- [21] Christop, I. (2024). nEMO: Dataset of Emotional Speech in Polish. *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, 12111-12116.
- [22] Burkhardt, F., Paeschke, A., Rolfes, M., Sendlmeier, W. F., & Weiss, B. (2005). A database of German emotional speech. *Proceedings of Interspeech 2005*, 1517-1520.
- [23] McFee, B., et al. (2015). librosa: Audio and Music Analysis in Python. *Proceedings of the 14th Python in Science Conference*, 18-25.

- [24] Paszke, A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. *Advances in Neural Information Processing Systems*, 32.
- [25] Askari, M., Shahzad, A., Khan, A. F., Fuzail, M., Aslam, N., & Tariq, M. (2025). EFFEC-TIVE SPEECH EMOTION RECOGNITION USING R-CNN & BLSTM. *Kashf Journal of Multidisciplinary Research*, 2, 293-309. DOI: 10.71146/kjmr514.
- [26] Mishra, S., & Rizvi, S. (2025). AI-Driven Speech Emotion Detection: A Systematic Approach To Voice-Based Sentiment Analysis. *International Journal of Computer Applications*, 187, 43-48. DOI: 10.5120/ijca2025924877.

This document was compiled using L^AT_EX with pdfT_EX 3.141592653-2.6-1.40.25 (T_EX Live 2023).