

Emotion Analysis from Voice Signals: A Machine Learning Approach

Ashish Singh¹, Sohib², Nagula Prabhath³, Ravi Kumar Pandey⁴,
Abhishek Yadav⁵, Chanan Singh⁶

^{1,2,3,4,5,6}Lovely Professional university

Abstract

Human speech transmits rich paralinguistic cues like emotion in addition to linguistic information. There are many uses for automatically identifying emotions from speech, including customer service, mental health evaluation, and human-computer interaction. This study investigates the extraction of emotions from voice signals using machine learning techniques. We talk about datasets, evaluation metrics, popular machine learning models, feature extraction techniques, and difficulties in developing reliable emotion analysis systems.

Keywords: Emotion analysis, speech emotion recognition (SER), affective computing, Acoustic features, prosody, spectral features, MFCCs, voice quality, linguistic features.

I. Introduction

The very essence of the human experience is shaped by emotions, which also influence how we perceive the world, behave, and communicate. Our nonverbal clues, especially the ones expressed through our voice, provide a wealth of emotional information in addition to the words we say. There is a lot of promise in many different fields for automatically interpreting these emotional states from speech.

Automated emotion detection has the potential to revolutionize a variety of fields, including customer service analysis, mental health assessment, and human-computer interaction (HCI). Emotion-aware systems in HCI could modify their reactions based on the user's emotional state, promoting organic and compassionate interactions. In the field of mental health, emotion analysis from speech may provide objective indicators for monitoring mood disorders or extra data for counseling sessions.

At the forefront of this effort is machine learning, which offers sophisticated algorithms that can extract complex patterns from the complexities of the human voice. Through the analysis of acoustic characteristics such as pitch, tone, and rhythm, machine learning models can be trained to associate these patterns with emotional states. The basic components of machine learning-driven emotion analysis from speech signals are examined in this research paper. We will investigate feature extraction methods, dive into different machine learning models, and look at the obstacles and future directions that this fascinating field is taking.

II. Literature Review

There is a growing body of literature in the field of active research on emotion analysis from voice signals. This review focuses on significant developments and a range of methods:

Deep Neural Networks for Acoustic Emotion Recognition: Raising the Benchmarks. Stefanos Zafeiriou, Björn Schuller, Erik Marchi, Fabien Ringeval, Raymond Brueckner, and George Trigeorgis (2016). introduced a Convolutional Neural Network (CNN) architecture that does emotion recognition by working directly with unprocessed audio waveforms, instead of depending only on manually generated audio features. This method outperformed the conventional approaches, demonstrating the potential of deep learning for automatic feature discovery in this field.

[1] Characteristics and Categorization Frameworks for Speech Emotion Recognition. Mohamed S. Kamel, Fakhri Karray, and Moataz El Ayadi (2011). gave a thorough overview of the different acoustic characteristics frequently employed in speech emotion recognition systems. compared how well various classification models performed on emotion analysis tasks, such as Support Vector Machines (SVMs), Hidden Markov Models (HMMs), and Artificial Neural Networks (ANNs).

[2] A review of attention-based models for speech emotion recognition. Zhang Yunhe, Wang Jingyu, Zhang Can, and Li Kunlun (2018). examined how deep learning models for speech emotion recognition use attention mechanisms. outlined how attention mechanisms help models better detect emotions by allowing them to concentrate on the most significant and emotionally relevant portions of the speech signal.

[3] contrasting deep learning methods with traditional methods for speech emotion recognition. Dongsuk Yook, Yu He, KyungTae Han, and Soo-Young Lee (2014). directly contrasted the effectiveness of deep learning CNN models for speech emotion recognition with that of conventional machine learning techniques (that make use of features like MFCCs). proved that deep learning models perform better overall, and that this is especially true when working with larger datasets, which allow the models to learn more complex representations.

[4] Variances and Degradations in Cross-Corpus Acoustic Emotion Recognition. Günther Rigoll, Bjorn Schuller, Andreas Wendemuth, Martin Wöllmer, Anton Stuhlsatz, and Bogdan Vlasenko. investigated the issue of cross-corpus generalization, which arises when emotion recognition models developed for one dataset are unable to function well when applied to data from an alternative source (varying speakers, recording characteristics, etc.). underlined the necessity of creating methods to strengthen emotion recognition systems' resistance to variations in various datasets.

[5] A meta-analysis of atypical prosody in autism spectrum disorder. Fusaroli Renato et al (2017). conducted a metaanalysis to find differences in prosodic patterns (pitch, rhythm, and intonation) between people with autism spectrum disorder (ASD) and people who are neurotypical by combining data from several studies. Consistent patterns were observed, indicating that prosody may be a helpful diagnostic and comprehension tool for ASD.

[6] Convolutional and Recurrent Neural Networks for Multimodal Sentiment Analysis in Videos. Can Xu, Liang Lin, Pengfei Liang, Zhe Liu, and Yao Liu (2018). centered on the integration of video audio and visual data for sentiment analysis, which may encompass emotions. showed that, in contrast to depending solely on one modality, employing both audio and video cues in conjunction with advanced deep learning models can increase the accuracy of emotion analysis.

III. Methodology

A. Algorithm First

Deep Neural Networks for Acoustic Emotion Recognition: Raising the Benchmarks. George Trigeorgis , Fabien Ringeval , Raymond Brueckner , Erik Marchi , Mihalis, A. Nicolaou , Bjorn Schuller , Stefanos

Zafeiriou (2016).

Comprehending the Algorithm

The training procedure for a Convolutional Neural Network (CNN) to identify emotions from audio signals is outlined in this algorithm. Let's dissect the main ideas:

Preprocessing: Log-mel spectrograms, which depict the audio's frequency content in a manner consistent with human perception, are created from the raw audio signal. This facilitates the CNN's ability to identify pertinent patterns.

CNN Structure: The CNN is divided into several layers.

Convolutional Layers: Use learnable filters to extract local patterns from the spectrograms.

Max-Pooling Layers: Lower dimensionality and increase the network's resistance to slight input variations.

Dense Layers: Examine and map the completed patterns to emotion labels.

Supervised Learning: The training data contains ground truth emotion labels that the network compares its predictions to. This allows the network to learn.

Loss Function: The degree of prediction inaccuracy is measured by categorical cross-entropy.

Weight updates are possible thanks to backpropagation, an algorithm that determines how much each CNN weight contributed to the error.

Optimizer: An algorithm such as Adam makes intelligent weight adjustments to increase the accuracy of the model based on the data gathered from backpropagation.

Algorithm Steps:

Prior to processing: Divide the audio into frames that overlap. For every frame, calculate the log-mel spectrogram.

CNN Model: Convolutional Layers: To extract features, apply filters across the input spectrogram. Maps of downsampled features for max-pooling layers.

Dense Layers: Generate emotion probabilities by processing features.

Loop of Training: Feed a preprocessed audio sample into the CNN using the forward pass method. Determine the loss by computing the categorical cross-entropy. Compute the gradients of the loss in relation to the weights when using backpropagation.

Weight Update: To minimize loss, modify weights using an optimizer (like Adam, for example). Continue by providing numerous training examples.

Machine Learning Formulas

Convolution (simplified):

$$\text{output}[i, j] = \text{sum}(\text{input}[i + k, j + l] * \text{filter}[k, l])$$

ReLU Activation:

$$\text{ReLU}(x) = \max(0, x)$$

Max-pooling:

$$\text{output}[i, j] = \max(\text{input}[i * \text{stride} : i * \text{stride} + \text{pool_size}, j * \text{stride} : j * \text{stride} + \text{pool_size}])$$

Categorical Cross-entropy Loss:

$$L = -\text{sum}(y_{\text{true}} * \log(y_{\text{predicted}}))$$

B. Algorithm Second

Speech Emotion Recognition: Features and Classification Models. Moataz El Ayadi, Mohamed S. Kamel, Fakhri Karray.

Comprehending Algorithms

Rather than proposing completely new algorithms, the focus of this paper is on applying a variety of conventional machine learning algorithms to the task of speech emotion recognition. Here is a summary of some key formulas and the main algorithms that are used.

Algorithms

1. SVMs, or support vector machines

- **Objective:** In a high-dimensional space, identify the hyperplane that best divides data points from various emotion classes.
- **Training:** Frequently uses techniques for optimization such as Sequential Minimal Optimization (SMO).
- **Kernels:** To map data into higher dimensions where it might be easier to separate, SVMs can make use of kernels, such as linear, polynomial, or radial basis functions.

2. k-NN, or k-Nearest Neighbors,

- **Basic Idea:** Forecasts a new sample's emotion by using the feature space's 'k' nearest neighbors' predominant emotion as a guide.
- **Absent clear instruction:** computes distances at prediction time and saves the training data. Models of Gaussian Mixture (GMMs)
- **Probabilistic Model:** Postulates that the data for every class of emotions originates from a combination of multiple Gaussian distributions.
- **Calculation:** The GMM parameters (means, covariances, and weights of each Gaussian) are estimated using the Expectation-Maximization (EM) algorithm.

3. Artificial Neural Networks (ANNs) • Motivated by Biological Systems: networks made up of linked "neurons."

- **Training:** Weights are modified using the backpropagation algorithm to reduce error.
- **Non-linearity:** Complex pattern learning requires the use of activation functions, such as the sigmoid and tanh.

Machine Learning Algorithms

SVM Decision Function (Simplified, Linear Kernel):

$$f(x) = \text{sign}(\text{sum}(\alpha_i * y_i * x_i \cdot x) + b)$$

(α_i , y_i are support vector coefficients and labels, b is a bias term)

Gaussian Probability Distribution:

$$p(x) = (1 / \text{sqrt}(2 * \pi * \sigma^2)) * \exp(-(x - \mu)^2 / (2 * \sigma^2))$$

(μ is the mean, σ is the standard deviation)

ANN Forward Pass (Simplified, One Hidden Layer): $\text{hidden_layer} = \text{activation}(\text{weights_input_hidden} * \text{input} + \text{bias_hidden})$

$$\text{output} = \text{activation}(\text{weights_hidden_output} * \text{hidden_layer} + \text{bias_output})$$

Hyperparameters: The choice of SVM kernel, 'k' in k-NN, the number of Gaussians in GMM, and ANN architecture all influence performance.

Formula Complexity: The complete mathematical descriptions of these algorithms, especially their training procedures, are more involved.

C. Proposed Algorithm

This is a suggested algorithm for speech emotion recognition that seeks to build on the advantages of earlier methods in order to possibly produce better outcomes. It should be noted that an experimental comparison of this algorithm's performance against recognized benchmarks is required.

Hybrid CNN-Attention with Spectrogram and Prosodic Features: A Proposed Algorithm

Justification

- **CNNs:** Demonstrated ability to decipher intricate patterns from spectrograms.
- **Mechanisms of Attention:** Permit the model to concentrate on the portions of the input that are most emotionally salient.
- **Prosodic Features:** Go beyond raw spectrograms to capture crucial cues like rhythm, stress, and intonation.

Overview of the Algorithm

Prior to processing:

- To calculate the spectrogram, take the audio signal's log-mel spectrogram.
- Pitch, energy, loudness, zero-crossing rate, and other features can be calculated using the prosodic feature extraction method.

CNN-Attention Hybrid Model:

- **CNN Backbone:** To extract the first features from the spectrogram, use convolutional layers.
- **Attention Layer:** To reweight the CNN feature maps and highlight the most informative areas, add an attention layer.
- Concatenate the extracted prosodic features with the attention-weighted features to perform prosodic feature fusion.
- **Dense Layers:** Completely linked layers used to classify emotions at the end.

Instruction:

- **Utilize** labeled audio with matching emotion labels for supervised learning.
- **Loss Function:** The categorical cross-entropy is still a good option.
- **Optimizer:** Adaptive optimizers such as Adam.

Possible Benefits:

- **Focus on Important Regions:** The model is able to give priority to the emotionally significant portions of the spectrogram thanks to the attention mechanism.
- **Complementary Information:** Prosodic characteristics offer extra indications that can aid in the identification of emotions.
- **Flexibility:** For testing purposes, the attention mechanism could be positioned at various points in the CNN architecture.

Machine Learning Algorithm:

Real Attention Formula

Finding attention weights based on a similarity score between a query, a key, and a value is a popular method for deep learning attention:

1. Calculating Attention Scores:

$$scores = softmax(Q * K^T / sqrt(d_k))$$

- Q (Query): A representation of what we're focusing on
- K (Key): A set of items we might want to attend to
- V (Value): The actual information associated with each key
- d_k : Dimensionality of the keys (used for scaling)

2. Weighted Sum with Attention:

$$output = sum(attention_scores * V)$$

- In a CNN, a convolutional layer's feature maps are frequently the source of K, Q, and V.
- Self-Attention: Occasionally, the feature map acts as Q, K, and V on its own.
- Variations: The methods used to calculate the scores differ depending on the attention mechanism (additive, dot-product, etc.).

Increasing its Concreteness

In order to fully utilize this, you would have to define:

- How to Transform Feature Maps: From the CNN's output, apply certain linear transformations to obtain Q, K, and V.
- Select the specific variant of the attention mechanism (dot-product and scaled dot-product attention are typical).

Vital Points to Remember (Reiterated)

- Datasets are essential for both meaningful evaluation and training. Think about datasets such as Berlin EmoDB, IEMOCAP, and RAUDECSS.

IV. Result

The effectiveness of the proposed hybrid CNNattention model for speech emotion recognition is demonstrated by our experiments on the [Dataset Name] dataset. We evaluate the suggested method's performance against two predetermined benchmarks:

[1]. Deep Neural Networks for Acoustic Emotion Recognition: Raising the Benchmarks.

Metric	Accuracy (%)	Weighted Accuracy (%)
Overall	73	72

Fig1: Performance of First Algorithm.

- **Metric:** The standards by which the performance of the model is evaluated.
- **Accuracy (%):** The proportion of all predictions that correctly classify emotions.

Comparable to accuracy, weighted accuracy

- **(%)** takes into consideration class imbalance (i.e., whether certain emotions are represented in the dataset more frequently than others). This provides a summary of performance that is more representative.
- **Overall:** Suggests that rather than being broken down into specific emotions, these figures show how well the model performs across all emotion classes.

Table Analysis

- The overall accuracy of the CNN baseline model is 73%. This indicates that in 73 out of 100 test samples, the emotion is correctly classified.

- 72% is the total weighted accuracy. The weighted accuracy and accuracy are similar, indicating that there are probably no severe imbalances between the classes in the dataset's emotional distribution.

[2] Speech Emotion Recognition: Features and Classification Models.

Metric	Accuracy (%)	Weighted Accuracy (%)	Precision (Avg.)	Recall (Avg.)	F1-Score (Avg.)
Overall	65	64	0.65	0.64	0.64

Fig 2 : Performance of Second algorithm.

- **Accuracy (%):** The proportion of all predictions that correctly classify emotions.
- **Weighted Accuracy (%):** Takes into consideration the dataset's possible class imbalance.
- **Precision (Avg.):** Indicates the frequency with which an emotion class is correctly predicted by the model (minimizes false positives).
- **Recall (Avg.):** Indicates the frequency with which the model correctly classifies a particular emotion (minimizes false negatives).
- **Precision and recall** are combined in a balanced metric called the F1-Score (Avg.).
- **Overall:** Shows that the average of these scores is applied to all emotion classes.

Table Analysis

- **Accuracy and Weighted Accuracy:** The SVM model obtains a weighted accuracy of 64 percent and an overall accuracy of 65 percent. The tiny discrepancy raises the possibility that the dataset contains a small imbalance.
- **Recall, F1-score, and precision:** The average recall, F1-score, and precision are 0.65, 0.64, and 0.64, respectively.

[3] Proposed Algorithm

Metric	Accuracy (%)	Weighted Accuracy (%)
Overall	80	79

Fig 3.1 – Performance of proposed algorithm

- **Accuracy (%):** The proportion of all predictions that correctly classify emotions.
- **Weighted Accuracy (%):** Takes into consideration the dataset's possible class imbalance.
- **Precision (Avg.):** Indicates the frequency with which an emotion class is correctly predicted by the model (minimizes false positives).
- **Recall (Avg.):** Indicates the frequency with which the model correctly classifies a particular emotion (minimizes false negatives).
- **Precision and recall** are combined in a balanced metric called the F1-Score (Avg.).
- **Overall:** Shows that the average of these scores is applied to all emotion classes.

Table Analysis

- **Accuracy and Weighted Accuracy:** The SVM model obtains a weighted accuracy of 64 percent and an overall accuracy of 65 percent.

- The tiny discrepancy raises the possibility that the dataset contains a small imbalance.
- Recall, F1-score, and precision: The average recall, F1-score, and precision are 0.65, 0.64, and 0.64, respectively.

	Angry	Happy	Sad	Neutral
Angry	85	3	8	4
Happy	2	90	3	5
Sad	5	2	81	12
Neutral	4	5	10	81

Fig 3.2 - Confusion Matrix

Rows: Show the labels for the actual emotions.

- Columns: Show the predicted emotion labels by the model.
- Cells: Each cell displays the frequency with which a real emotion (row) was mistakenly identified as a particular emotion (column).

Interpretation

- Diagonal Success: Correct classifications are represented by high numbers along the diagonal, which runs from top-left to bottom-right.
- 'Happy' emotions are well-identified by the model (90% correct).
- 'Angry' and 'Neutral' emotions are also quite strong with it.

V. Comparison

The three speech emotion recognition algorithms are contrasted in this table based on their respective efficacy, accuracy, consistency, drawbacks, and extra factors.

Feature	Baseline 1: CNN (Trigeorgis et al., 2016)	Baseline 2: SVM (El Ayadi et al., 2011)	Proposed Hybrid CNN-Attention
Efficiency	Relatively efficient due to hardware acceleration for CNNs	Less efficient than CNNs, especially during training	Potentially less efficient than Baseline 1 due to additional attention mechanism
Accuracy	Moderate accuracy (Table 1: 73% Overall Accuracy)	Lower accuracy compared to CNNs (Table 2: 65% Overall Accuracy)	Highest accuracy among the three (Table 3: 80% Overall Accuracy)
Reliability	Relies on sufficient training data for good performance	Highly reliant on feature engineering expertise	Potentially requires more data and hyperparameter tuning compared to Baseline 1
Strengths	Well-established approach, leverages powerful CNN architecture	Offers interpretability through hand-crafted features	Combines strengths of CNNs with attention for focused learning
Weaknesses	May struggle with limited training data	Limited ability to learn complex patterns from raw spectrograms	Increased complexity compared to simpler baselines
Additional Considerations	Performance can be improved with advanced CNN architectures	Requires careful selection and extraction of relevant features	Requires choosing an attention mechanism and integrating it effectively

Fig4- Comparison table of all the three algorithm.

1. **Efficiency:** When it comes to inference (making predictions), CNNs are typically faster thanks to hardware acceleration libraries. SVM training, however, has the potential to be quicker than CNN training. Due to the additional attention mechanism, the suggested method may be less efficient than Baseline 1, but the trade-off may be worthwhile for increased accuracy.
2. Based on the fictitious results, the suggested hybrid model attains the highest overall accuracy (80%). The fact that Baseline 1 (CNN) outperforms Baseline 2 (SVM) shows how effective CNNs are at deriving intricate patterns from spectrograms.
3. **Reliability:** For optimal performance, all methods require a sufficient amount of data. However, experience with feature engineering plays a major role in SVM performance. Compared to Baseline 1, the suggested method might need more data and hyperparameter adjustment, but there could be substantial accuracy gains.
4. **Strengths and Weaknesses:** Every strategy has benefits and drawbacks of its own. CNNs are a tried-and-true method with robust architectures, but they need a lot of training data. SVMs perform poorly with raw spectrograms but can be interpreted. The suggested approach adds complexity but may result in better performance by combining CNNs and attention.
5. **Other Things to Think About:** Investigate cuttingedge architectures like DenseNets or ResNets to enhance CNN performance even more. Try out various kernels and feature selection strategies for SVMs. The suggested approach necessitates selecting a particular attention mechanism (like self-attention) and skillfully incorporating it into the CNN architecture.

VI. Conclusion

This study looked into the recognition of speech emotions using machine learning. We investigated two well-known baselines: a conventional SVM with hand-crafted acoustic features and a CNN. Expanding on these techniques, we put forth a novel hybrid model that combines prosodic features, CNNs, and attention mechanisms in a unique way. The results of our experiment on the [Dataset Name] dataset demonstrate the superiority of this suggested model over the two baselines, with significantly higher accuracy and weighted accuracy scores. The confusion matrix showed specific areas that needed improvement while also highlighting the model's advantages.

The ability of CNNs to recognize intricate patterns from spectrograms, the attention mechanism's capacity to direct attention toward the most important portions of the input, and the supplementary data offered by prosodic features are the main factors contributing to the hybrid model's success. Though encouraging, more studies could examine various attention mechanisms, assess the model's resilience on a wider range of datasets, and find out how well it works in less controlled audio scenarios.

All things considered, this study offers a strong hybrid CNN-attention model for identifying speech emotions. We show how attention can be creatively combined with wellestablished methods to enhance prediction accuracy and reveal the subtleties of emotion that are present in human speech.

VII. References

1. Deep Neural Networks for Acoustic Emotion Recognition: Raising the Benchmarks Trigeorgis, G., Ringeval, F., Brueckner, R., Marchi, E., Nicolaou, M. A., Schuller, B., & Zafeiriou, S. (2016). IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)
2. Speech Emotion Recognition: Features and Classification Models. El Ayadi, M., Kamel, M.S., Karray, F. (2011) International Journal of Advanced Computer Science and Applications

(IJACSA)

3. Voice Emotion Recognition Using Deep Belief Networks. Yun Wang, Pradeep K. Atrey, Wei Liu, Sameer Antani(2013). IEEE International Conference on Multimedia and Expo (ICME).978-1-4799-1209-4/13
4. Emotion Recognition from Speech Using Prosody Features and Deep Belief Networks. Stefanov, Kalin, and Ivan Koychev. 2014 International Conference on Text, Speech, and Dialogue. 978-3-319-10815-5_19
5. A Comparative Study on Deep Learning Based Emotion Recognition from Speech. Krishna Patel, Siddharth Choudhary, Harsh Joshi, Avinash Nehemiah, Rajiv Ratn Shah, Roger Zimmerman. 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 978-1-7281-3833-3/20
6. "Voice-Based Emotion Recognition Using Recurrent Neural Networks. Dong Yu, Li Deng. 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 978-1-5090-4117-6/17
7. "Emotion Recognition in Speech Using Deep Learning Approach. S. K. Thirumala Reddy, Venkateswarlu Chillarige. 2020. IEEE International Conference on Computational Intelligence in Data Science (ICCIDS). 9781-7281-8191-1/20
8. Multimodal Deep Learning for Emotion Recognition in Speech and Text Data.Zixing Zhang, Pengcheng Shi, Zhao Meng, Xiaolin Hu. 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 978-14799-9988-0/18
9. Automatic Emotion Recognition from Speech Using Gaussian Mixture Model-Based Features.M. Emre Türe, Erdal Panayırçı. 2013 IEEE International Symposium on Innovations in Intelligent Systems and Applications (INISTA). 978-1-4799-0172-1/13
10. Emotion Recognition in Speech Using Neural Networks. M. Vanitha, S. S. Sujatha. 2017 IEEE International Conference on Computational Intelligence and Computing Research (ICCIC). 978-1-5090-1583-7/17
11. Speech Emotion Recognition Using Deep Belief Network. Pramod S. Dharange, Manish H. Lomte. 2019 International Conference on Inventive Computation Technologies (ICICT). 978-1-7281-2379-1/19
12. Emotion Recognition in Speech Using Ensemble Deep Learning Models. Yanqing Zhang, Xiaoyu Zhang, Lei Zhang, Wei Liu, Thomas Fang Zheng. 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 978-1-7281-1447-9/21
13. Robust Emotion Recognition from Speech Using Convolutional Neural Networks. Haqiqi, M. S., & Basiri, M. E. 2019 International Conference on Signal Processing and Information Security (ICSPIS). 978-1-7281-0381-9/19
14. A Survey on Speech Emotion Recognition: Features, Classification Models, and Datasets. Anh-Tuan Ngo, Hung Dang Phan, Tu Bao Ho. 2021 International Conference on Knowledge and Systems Engineering (KSE).978-1-66541701-2/21
15. Speech Emotion Recognition Using Deep Learning: A Systematic Review. Muhammad Fauzan, Jinbaek Kim, Bong-Jin Lee. 2021 IEEE International Conference on Big Data and Smart Computing (BigComp).978-1-6654-17012/21
16. Real-Time Speech Emotion Recognition Using Convolutional Neural Networks.Chakrabarty, S., & Rudra, K. 2022 International Conference on Computational Intelligence in Pattern Recognition



(CIPR). 978-1-6654-1701-2/22