

SPEECH EMOTION RECOGNITION USING QUATERNION CONVOLUTIONAL NEURAL NETWORKS

Aneesh Muppidi and Martin Radfar,

Department of Computer Science, Stony Brook University, YN, USA

aneeshmuppidi19@gmail.com, radfar@cs.stonybrook.edu

Abstract

Although speech recognition has become a widespread technology, inferring emotion from speech signals still remains a challenge. To address this problem, this paper proposes a quaternion convolutional neural network (QCNN) based speech emotion recognition (SER) model in which Mel-spectrogram features of speech signals are encoded in an RGB quaternion domain. We show that our QCNN based SER model outperforms other real-valued methods in the Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS, 8-classes) dataset, achieving, to the best of our knowledge, state-of-the-art results. The QCNN also achieves comparable results with the state-of-the-art methods in the Interactive Emotional Dyadic Motion Capture (IEMOCAP 4-classes) and Berlin EMO-DB (7-classes) datasets. Specifically, the model achieves an accuracy of 77.87%, 70.46%, and 88.78% for the RAVDESS, IEMOCAP, and EMO-DB datasets, respectively. In addition, our results show that the quaternion unit structure is better able to encode internal dependencies to reduce its model size significantly compared to other methods.

Index Terms: Speech Emotion Recognition, Signal Processing, Quaternion Deep Learning, Convolutional Neural Networks

1. Introduction

Speech emotion recognition (SER) is an active, yet challenging area of research that has many important implications in technologies such as automated healthcare, clinical trials, voice assistants, psychological therapy, emergency responders, call centers, video games, robot-human interactions, and more. Nonetheless, recognizing emotion from speech signals is difficult due to many reasons, namely: qualitative properties of identifying emotion, background noise, variable individual-specific accentuation, weak representation of grammatical and semantic knowledge, temporal and spectral characteristics of the signal due to fast or slow speech, dynamical characteristics such as soft or loud voices, and multiple voices [1–3]. In addition, if an application can learn high-level features from speech for emotion classification, it is possible to replicate or imitate emotion from a text-to-speech context, further improving the aforementioned fields of impact [4].

In general, SER is a classification task in which we extract features from a set of emotion-labeled speech signals and use the pair of feature and label to train a classifier which is commonly a deep neural architecture. Features such as MFCC, chromograms, spectral contrast features, and Tonnetz representations have been used in previous neural based SER models [5–7]. Neural based SER models usually leverage n -dimensional convolutional neural networks (CNNs) [8], recurrent neural networks (RNNs) [9], Long short-term memory networks (LSTMs) [10], or as combinations and variations of these techniques [11].

Recently, the quaternion based neural networks have been shown to better encode features than the standard real-valued based neural networks [12]. The use of the quaternion domain has shown a unique ability in color image processing to capture all three channels of the RGB domain without the typical information loss of separating these channels, while reducing the size of its model in comparison to real-valued models [13].

In this paper, to the best of our knowledge, we propose the first quaternion based model that is able to encode an RGB representation of Mel-spectrograms for the application of SER classification. Using this approach, we propose a quaternion convolutional neural network (QCNN) to train on benchmark SER datasets in the quaternion domain with quaternion converted elements of a CNN. The results, conducted on three public SER datasets, show that the QCNN model is able to yield state-of-the-art results on RAVDESS and outperforms all but one competitive models in both IEMOCAP and EMO-DB in terms of accuracy. In addition, our proposed model has the smallest size compared to the previous competitive architectures.

The rest of this paper is organized as follows: In Section 2 we mention some related works. In Section 3, we give a brief overview of quaternion operations and describe the proposed model. In Section 4 we detail our experiments and results. Finally, we draw the conclusion and give the future directions in Section 5.

2. Relation to prior works

A majority of previous literature on SER methods consist of variations CNNs such as a 1d-CNN for waveform features, CNN-LSTMs, transfer CNN models [14]. In addition, many methods have also focused on RNNs [15, 16] to discover new features. In the context of quaternion models, [12] introduced the ability to generate quaternion-based dense layers. Extending this work, QCNNs were developed for color image processing [13]. The significance of quaternion models in the context of signal processing was originally shown by [17] for automatic speech recognition; however, this model encoded views of a time-frame frequency as quaternions. Our work extends previous literature in that we exploit the decibel visualization of Mel-spectrograms as RGB arrays to encode in the quaternion domain rather than encoding the waveform features as pure quaternions or using a real-valued model.

3. Model

3.1. Preprocessing

Our feature generation pipeline consists of two parts: Mel-spectrogram array generation and RGB conversion. After the wav files were split into emotion-labelled folders, the Fourier transform was computed on each of the speech waveforms

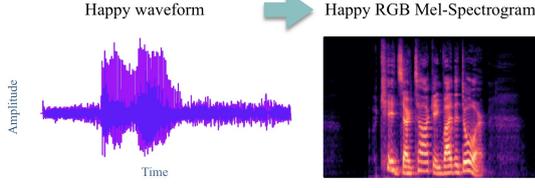


Figure 1: The happy waveform converted to an RGB Mel-Spectrogram

to transform from the time domain to the frequency domain. While the time axis scaled linearly by seconds, the frequency axis was non-linearly transformed using a Mel transform, which is defined as:

$$F_{\text{mel}} = 2,595 \log_{10} \left(1 + \frac{f}{700} \right) [\text{dB}] \quad (1)$$

where f denotes the frequency axis. The spectral amplitudes of the arrays were normalized between their dataset specific Max[dB] and Min [dB]. The mel-spectrogram was then converted to an image, with a color scale to represent the amplitudes (see example given in Figure 1). The RGB images were then split into training (80%) and test (20%) datasets.

3.2. Quaternion Convolutional Neural Networks

3.2.1. Quaternion Algebra

Briefly, a quaternion is hypercomplex number which is an extension to complex numbers in 3-dimensional space with additionally imaginary parts \mathbf{j} and \mathbf{k} . Thus, a quaternion Q in domain Q can be represented as $Q = r_0 + r_1\mathbf{i} + r_2\mathbf{j} + r_3\mathbf{k}$ where $r \in R$ and \mathbf{i} , \mathbf{j} , and \mathbf{k} are imaginary parts. Like complex number, Quaternions obey specific rules:

$$\mathbf{i}^2 = \mathbf{j}^2 = \mathbf{k}^2 = \mathbf{i}\mathbf{j}\mathbf{k} = -1 \quad (2)$$

and the addition is defined as

$$\begin{aligned} Qa + Qb &= (r_{a_0} + r_{b_0}) \\ &+ (r_{a_1} + r_{b_1})\mathbf{i} \\ &+ (r_{a_2} + r_{b_2})\mathbf{j} \\ &+ (r_{a_3} + r_{b_3})\mathbf{k}. \end{aligned} \quad (3)$$

Scalar multiplication is defined as

$$xQ = xr_0 + xr_1\mathbf{i} + xr_2\mathbf{j} + xr_3\mathbf{k}. \quad (4)$$

Element multiplication is defined as:

$$\begin{aligned} Q_a \times Q_b &= (r_{a_0}r_{b_0} - r_{a_1}r_{b_1} - r_{a_2}r_{b_2} - r_{a_3}r_{b_3}) \\ &+ (r_{a_0}r_{b_1} + r_{a_1}r_{b_0} + r_{a_2}r_{b_3} - r_{a_3}r_{b_2})\mathbf{i} \\ &+ (r_{a_0}r_{b_2} - r_{a_1}r_{b_1} + r_{a_2}r_{b_0} + r_{a_3}r_{b_1})\mathbf{j} \\ &+ (r_{a_0}r_{b_3} + r_{a_1}r_{b_2} - r_{a_2}r_{b_1} + r_{a_3}r_{b_0})\mathbf{k}, \end{aligned} \quad (5)$$

and rotation along axis p is defined as

$$\hat{R} = p\hat{Q}p^{-1} \quad (6)$$

where

$$p = \cos\left[\frac{\theta}{2}\right] + (p_1 \times \mathbf{i} + p_2 \times \mathbf{j} + p_3 \times \mathbf{k}) \sin\left[\frac{\theta}{2}\right], \quad (7)$$

$$\hat{Q} = r_{q_0} + r_{q_1} \times \mathbf{i} + r_{q_2} \times \mathbf{j} + r_{q_3} \times \mathbf{k}, \quad (8)$$

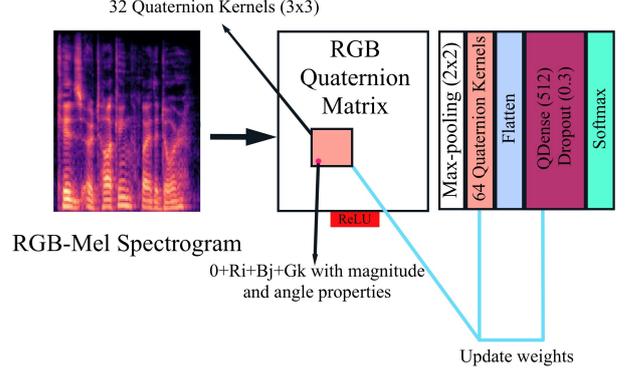


Figure 2: A high-level block diagram of the proposed QCNN SER system.

3.2.2. Quaternion Convolution

We implement a quaternion color kernel similar to the one proposed in [13]. We do this by defining the 50×75 colored Mel-spectrogram image in a 3D vector color space as a quaternion matrix $\hat{C} = [\hat{c}_{nn'}] \in H^{50 \times 75}$. We can thus represent the color channels in C as:

$$\hat{C} = 0 + Ri + Gj + Bk \quad (9)$$

where R, G , and $B \in R^{50 \times 75}$ are the red, green, and blue channels, respectively. By representing a pixel in 3D vector color space, the proposed model also deploys a quaternion kernel convolving around the input, in which the element $\hat{w}_{ll'}$ of the kernel \hat{W} can be defined as:

$$\hat{w}_{ll'} = s_{ll'} \left(\cos \frac{\theta_{ll'}}{2} + \left(\frac{\sqrt{3}}{3} (i + j + k) \right) \sin \frac{\theta_{ll'}}{2} \right) \quad (10)$$

where $\theta_{ll'} \in [-\pi, \pi]$ and $s_{ll'} \in R$. Thus the quaternion convolution operation $*$ can be expressed as:

$$\hat{C} \times \hat{W} = \hat{F} = [\hat{f}_{kk'}] \in H^{((50-L+1) \times (75-L+1))} \quad (11)$$

where \hat{W} is a $L \times L$ quaternion kernel and $[\hat{f}_{kk'}]$ is defined as:

$$\hat{f}_{kk'} = \sum_{l=1}^L \sum_{l'=1}^L \frac{1}{s_{ll'}} \hat{w}_{ll'} \hat{c}_{(k+l)(k+l')} \hat{w}_{ll'}. \quad (12)$$

By implementing quaternion convolution kernels to perform rotation and scaling operations, the model finds better representative features in the decibel color waveform space than real-valued convolution kernels. Additionally, the real-valued kernels apply scaling and pixel transformation separately to the three axes of color and create single-channel feature maps; whereas the quaternion kernels can more intuitively capture the color space as a whole without any type of color information loss since the channels are interrelated as a 3D vector.

3.2.3. Connecting Layers

The quaternion convolution layer is connected to other typical layers to construct the fully connected neural network. Specifically, max-pooling is performed on the imaginary parts and ReLU is used to reset invalid quaternion vector rotations to the nearest point in color space. In addition, the imaginary parts of

the output of a quaternion layer are represented as 3 real numbers, such that a real-valued and vectorized output can be obtained to be connected to real-valued fully-connected layer. A softmax layer is then connected to this real-valued layer to train the QCNN model.

3.2.4. Model Architecture

Figure 2 shows a block diagram of the proposed model. The RGB-Mel spectrogram is transformed into a 3D color vector space with rotation magnitudes for the initial input layer. This layer is then followed by the quaternion convolutional layer, which consists of 32 quaternion kernels (all quaternion kernels, proposed in this paper, cover a 3x3 2D coordinate space) that perform rotation and scaling operations. ReLU is used to reset invalid quaternion vector rotations to the nearest point in color space. Afterward, max-pooling is performed on the imaginary color properties. This is followed by 64 quaternion kernels. The layer is then flattened to vectorize the quaternion layer to a real-valued layer. The real-valued dense layer is used to train the model by computing the gradient of the categorical cross-entropy loss function (which operates on softmax outputs for the classes of emotion).

3.3. Training

Weights are initialized as an imaginary quaternion with a uniform distribution in the interval $[0, 1]$. The imaginary unit is then normalized, with well-known criteria described by [19] to generate the quaternion weight \hat{W} which was already defined above. Backpropagation is used to update the quaternion weights by applying the quaternion chain rule. Additionally, the loss function can be computed on the real-valued layer; specifically, the categorical cross-entropy function is used. Furthermore, the model used Adam as its stochastic optimizer and overfitting was addressed with the use of dropout (with a probability of 0.3). The model was trained for 50 epochs with 9 steps per validation epoch. The full model was implemented in Python with Tensorflow 2.3.0 as the deep learning backend.

4. Experiments

In order to measure the accuracy and effectivity of the model, QCNN is tested and compared against the most competitive models in SER on the RAVDESS, IEMOCAP, and EMO-DB datasets. In addition, we also compared the size of the models.

4.1. Speech Datasets

4.1.1. RAVDESS

The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) [20] consists of speech samples of 8 categorical emotions— calm, happy, sad, angry, fearful, surprise, and disgust —spoken by 24 professional actors, vocalizing lexically-matched statements in a neutral North American accent. The classification accuracy of QCNN for the 8 classes of emotion on RAVDESS was compared with results from CNN-LSTM [21], a pre-feature extracted SVM [22], GResNets [23], a controlled human accuracy study [RAVDESS], and Deep-CNN [14]. The results are summarized in table I.

Table 1: Classifier performance on RAVDESS

Model	Accuracy (%)
CNN-LSTM [21]	53.08
Feature SVM	60.10
GResNets [23]	65.97
Human Accuracy [20]	67.00
Deep-CNN [14]	71.67
QCNN [ours]	77.87

4.1.2. IEMOCAP

The Interactive Emotional Dyadic Motion Capture (IEMOCAP) [24] database is an acted, multimodal and multispeaker database. It consists of 12 hours of audiovisual data, including video, speech, motion capture of face, text transcriptions. This paper makes use of the speech samples of 4 categorical emotions— anger, happiness, sadness, and neutrality —spoken by 10 (5 female and 5 male) professional actors in fluent English. The classification accuracy of QCNN for the 4 classes of emotion on IEMOCAP was compared with results from CNN-LSTM [21], DialogueRNN [25], BiF-AGRU [26], Deep-CNN [14], DialogueGCN [27], BiERU-lc [28], and DSCNN [29]. The results are summarized in table II.

4.1.3. EMO-DB

The Berlin EMO-DB [30] consists of speech samples representing 7 categorical emotions —anger, happiness, sadness, fear, disgust, boredom, and neutral speech—spoken by 10 professional actors (5 female and 5 male) in German. Each speaker acted 7 emotions for 10 different utterances (5 short and 5 long) with emotionally neutral linguistic contents. The classification accuracy for the 7 classes of emotion on EMO-DB was compared for the Multi-task LSTM [31], RNN [21], CNN-LSTM [21], and Deep-CNN [14] methods with QCNN. The results are summarized in table III.

4.2. Model Size

The only model sizes that were reported for above models are Deep-CNN and CNN-LSTM. However, both being competitive models, they can be used as a bench mark to compare the QCNN model size. The results are summarized in Table IV.

4.3. Discussion

The QCNN model presented yields state-of-the-art results on RAVDESS by 6%. Although the model does not achieve state-of-the-art results on EMO-DB or IEMOCAP, in both cases it underperforms the most competitive model by only $<2\%$. The balance of accuracy over model size is also considered in this paper, and although the model slightly underperforms in the IEMOCAP and EMO-DB datasets, the size of the model is significantly less. This is attributed to the ability to store RGB arrays in a quaternion domain to reduce model size [12]. However, in this paper, we did not explore encoding other features of waveforms such as Mel-frequency cepstral coefficients (MFCC), chromograms, spectral contrast features, or Tonnetz representations as pure quaternions. It is possible that separate Quaternion models encoding these features, separately or individually, could outperform the model presented in this paper. Additionally, it is possible that an auxiliary neural network can be used in conjunction with the QCNN to capture higher-

Table 2: Classifier performance on IEMOCAP

Model	Unweighted Accuracy (%)
CNN-LSTM [21]	50.17
DialogueRNN [25]	63.40
BiF-AGRU [26]	63.50
Deep-CNN [14]	64.30
DialogueGCN [27]	65.25
BiERU-1c [28]	66.11
QCNN [ours]	70.46
DSCNN DSCNN[29]	72.00

Table 3: Classifier performance on EMO-DB

Model	Unweighted Accuracy (%)
Multi-task LSTM [31]	58.14
RNN[21]	63.21
CNN-LSTM [21]	69.72
QCNN [ours]	88.78
Deep-CNN [14]	90.01

level features. General performance can also be improved by using effective waveform data augmentation techniques. This research can also be extended to real-time SER for quaternion blocks to compute and process on speech-waveforms. However, for such a model it is expected to take advantage of other features rather than an RGB domain of Mel-spectrograms because of delta-computational expense.

5. Conclusion

Current SER research is a dynamic, complex task relying on feature extraction and computational classification. To capture this highly complex qualitative information from speech waveforms, multiple feature extraction and classification techniques have been introduced in literature. However, many machine learning-based methods focus on higher-level features in the real-valued space. This paper reports a unique approach to feature and network encoding by using a quaternion structural model. Specifically, we encode the RGB domain of Mel-spectrogram features in a quaternion input and use custom quaternion convolutional layers to learn features in a quaternion space. These layers are implemented in a standard neural network structure to train on benchmark datasets such as RAVDESS, IEMOCAP, and EMO-DB. QCNN is reported to yield an accuracy of 77.87%, 70.46%, and 88.78% for the RAVDESS, IEMOCAP, and EMO-DB datasets, respectfully. In comparison to other competitive methods, QCNN achieves state-of-the-art results on RAVDESS, and underperforms only one method in both IEMOCAP and EMO-DB. Additionally, QCNN is able to exploit its quaternion encoding for a reduced model size, confirming previous literature on the topic as well as providing an opportunity for deployment on lightweight machines such as voice-assistant devices. Given the significant results in both accuracy and computational complexity, QCNN is expected to be a foundational block for the continuation of SER research.

Table 4: Model sizes (Mb)

Model	RAVDESS	EMO-DB	IEMOCAP
Deep CNN	74.5	69.5	88.5
CNN-LSTM	111.1	94.4	128.3
QCNN	42	31.2	67.7

6. References

- [1] D. Rethage, J. Pons, and X. Serra, "A wavenet for speech denoising," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5069–5073.
- [2] M. Benzeghiba, R. De Mori, O. Derou, S. Dupont, T. Erbes, D. Jouvet, L. Fissore, P. Laface, A. Mertins, C. Ris *et al.*, "Automatic speech recognition and speech variability: A review," *Speech communication*, vol. 49, no. 10-11, pp. 763–786, 2007.
- [3] S. Ramakrishnan and I. M. El Emary, "Speech emotion recognition approaches in human computer interaction," *Telecommunication Systems*, vol. 52, no. 3, pp. 1467–1478, 2013.
- [4] A. Ashar, M. S. Bhatti, and U. Mushtaq, "Speaker identification using a hybrid cnn-mfcc approach," in *2020 International Conference on Emerging Trends in Smart Technologies (ICETST)*. IEEE, 2020, pp. 1–4.
- [5] S. An, Z. Ling, and L. Dai, "Emotional statistical parametric speech synthesis using lstm-rnns," in *2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2017, pp. 1613–1616.
- [6] S. Pascual, A. Bonafonte, and J. Serra, "Segan: Speech enhancement generative adversarial network," *arXiv preprint arXiv:1703.09452*, 2017.
- [7] E. J. Humphrey, T. Cho, and J. P. Bello, "Learning a robust tonnetz-space transform for automatic chord recognition," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012, pp. 453–456.
- [8] D. Palaz, M. M. Doss, and R. Collobert, "Convolutional neural networks-based continuous speech recognition using raw speech signal," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 4295–4299.
- [9] Y. Miao, M. Gowayyed, and F. Metze, "Eesen: End-to-end speech recognition using deep rnn models and wfst-based decoding," in *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*. IEEE, 2015, pp. 167–174.
- [10] A. Graves, N. Jaitly, and A.-r. Mohamed, "Hybrid speech recognition with deep bidirectional lstm," in *2013 IEEE workshop on automatic speech recognition and understanding*. IEEE, 2013, pp. 273–278.
- [11] M. Wollmer, B. Schuller, F. Eyben, and G. Rigoll, "Combining long short-term memory and dynamic bayesian networks for incremental emotion-sensitive artificial listening," *IEEE Journal of selected topics in signal processing*, vol. 4, no. 5, pp. 867–881, 2010.
- [12] T. Isokawa, T. Kusakabe, N. Matsui, and F. Peper, "Quaternion neural network and its application," in *International conference on knowledge-based and intelligent information and engineering systems*. Springer, 2003, pp. 318–324.
- [13] X. Zhu, Y. Xu, H. Xu, and C. Chen, "Quaternion convolutional neural networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 631–647.
- [14] D. Issa, M. F. Demirci, and A. Yazici, "Speech emotion recognition with deep convolutional neural networks," *Biomedical Signal Processing and Control*, vol. 59, p. 101894, 2020.
- [15] W. Lim, D. Jang, and T. Lee, "Speech emotion recognition using convolutional and recurrent neural networks," in *2016 Asia-Pacific signal and information processing association annual summit and conference (APSIPA)*. IEEE, 2016, pp. 1–4.

- [16] A. Satt, S. Rozenberg, and R. Hoory, "Efficient emotion recognition from speech using deep learning on spectrograms." in *Interspeech*, 2017, pp. 1089–1093.
- [17] T. Parcollet, Y. Zhang, M. Morchid, C. Trabelsi, G. Linares, R. De Mori, and Y. Bengio, "Quaternion convolutional neural networks for end-to-end automatic speech recognition," *arXiv preprint arXiv:1806.07789*, 2018.
- [18] E. O. Brigham and R. Morrow, "The fast fourier transform," *IEEE spectrum*, vol. 4, no. 12, pp. 63–70, 1967.
- [19] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 2010, pp. 249–256.
- [20] S. R. Livingstone and F. A. Russo, "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.
- [21] J. Parry, D. Palaz, G. Clarke, P. Lecomte, R. Mead, M. Berger, and G. Hofer, "Analysis of deep learning architectures for cross-corpus speech emotion recognition." in *INTERSPEECH*, 2019, pp. 1656–1660.
- [22] P. Shegokar and P. Sircar, "Continuous wavelet transform based speech emotion recognition," in *2016 10th International Conference on Signal Processing and Communication Systems (ICSPCS)*. IEEE, 2016, pp. 1–8.
- [23] Y. Zeng, H. Mao, D. Peng, and Z. Yi, "Spectrogram based multi-task audio classification," *Multimedia Tools and Applications*, vol. 78, no. 3, pp. 3705–3722, 2019.
- [24] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan, "Iemocap: Interactive emotional dyadic motion capture database," *Language resources and evaluation*, vol. 42, no. 4, p. 335, 2008.
- [25] N. Majumder, S. Poria, D. Hazarika, R. Mihalcea, A. Gelbukh, and E. Cambria, "Dialoguernn: An attentive rnn for emotion detection in conversations," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 6818–6825.
- [26] C. Li, J. Jiao, Y. Zhao, and Z. Zhao, "Combining gated convolutional networks and self-attention mechanism for speech emotion recognition," in *2019 8th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW)*. IEEE, 2019, pp. 105–109.
- [27] D. Ghosal, N. Majumder, S. Poria, N. Chhaya, and A. Gelbukh, "Dialoguecn: A graph convolutional neural network for emotion recognition in conversation," *arXiv preprint arXiv:1908.11540*, 2019.
- [28] Y. Li, T. Zhao, and T. Kawahara, "Improved end-to-end speech emotion recognition using self attention mechanism and multitask learning." in *Interspeech*, 2019, pp. 2803–2807.
- [29] S. Kwon *et al.*, "A cnn-assisted enhanced audio signal processing for speech emotion recognition," *Sensors*, vol. 20, no. 1, p. 183, 2020.
- [30] F. Burkhardt, A. Paeschke, M. Rolfes, W. F. Sendlmeier, and B. Weiss, "A database of german emotional speech," in *Ninth European Conference on Speech Communication and Technology*, 2005.
- [31] S. Goel and H. Beigi, "Cross lingual cross corpus speech emotion recognition," *arXiv preprint arXiv:2003.07996*, 2020.