

Effect of Attention and Self-Supervised Speech Embeddings on Non-Semantic Speech Tasks

Payal Mohapatra*
Northwestern University
Evanston, Illinois, USA
payal.mohapatra@northwestern.edu

Akash Pandey*
Northwestern University
Evanston, Illinois, USA
akash.pandey@northwestern.edu

Yueyuan Sui*
Northwestern University
Evanston, Illinois, USA
yueyuansui2024@u.northwestern.edu

Qi Zhu
Northwestern University
Evanston, Illinois, USA
qzhu@northwestern.edu

ABSTRACT

Human emotion understanding is pivotal in making conversational technology mainstream. We view speech emotion understanding as a perception task which is a more realistic setting. With varying contexts (languages, demographics etc.) different *share* of people perceive the same speech segment as a non-unanimous emotion. As part of the ACM Multimedia 2023 Computational Paralinguistics Challenge (ComParE) in the **EMotion Share** track, we leverage their rich dataset of multilingual speakers and multi-label regression target of 'emotion share' or perception of that emotion. We demonstrate that the training scheme of different foundation models dictates their effectiveness for tasks beyond speech recognition, especially for non-semantic speech tasks like emotion understanding. This is a very complex task due to multilingual speakers, variability in the target labels, and inherent imbalance in the regression dataset. Our results show that HuBERT-Large with a self-attention-based light-weight sequence model provides 4.6% improvement over the reported baseline.

CCS CONCEPTS

• **Computing methodologies** → **Natural language processing.**

KEYWORDS

Emotion share, large-language model, attention

ACM Reference Format:

Payal Mohapatra, Akash Pandey, Yueyuan Sui, and Qi Zhu. 2023. Effect of Attention and Self-Supervised Speech Embeddings on Non-Semantic Speech Tasks. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3581783.3612855>

*All authors contributed equally to this research. Payal worked on the large language model embeddings for feature extraction, Akash worked on the attention module, and Yueyuan worked on the data preparation and model selection.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0108-5/23/10...\$15.00
<https://doi.org/10.1145/3581783.3612855>

1 INTRODUCTION

Human emotion understanding is vital to engage in meaningful social interactions. With the rising use of voice-assisted technology and the use of natural language as a key Human-Machine interface, it is important that we develop techniques to understand the affective aspects of speech [19]. We need to correctly identify and characterize the tone of speech and not merely understand its semantics of 'what' has been said. It is crucial to understand 'how' it has been said. Classifying a speech segment as one type of emotion is oversimplifying the way humans perceive language [6]. On many occasions under various contexts, different groups of people perceive the same emotion differently. This gives rise to a viewpoint that we can assign a *perceivability share* for an emotion in speech.

Traditionally, speech features like mel-frequency cepstral coefficients (MFCCs) [11, 12], filterbanks, fundamental frequency, energy, zero-crossing rate, chroma-based features, and their feature functionals [1, 10] are used with a rule-based or neural network classifier to recognize the emotions in speech. In recent times, tremendous progress of deep-learning models in the field of Natural Language Processing has positively boosted the performance of paralinguistic speech tasks [4, 7, 9] like emotion recognition as well. Several works in the literature advocate the use of self-supervised pipelines (SSL) [14, 21, 22] for Automatic-Speech-Recognition tasks for SER. The two foundation models that emerge as the current state-of-the-art feature extractors for ASR tasks are HuBERT [8] and wav2vec2.0 [2].

Most if not all the previous works on understanding human emotions from speech can be categorized as Speech Emotion Recognition (SER). Such studies suffer from a few limitations - 1) the datasets consist of actors emulating certain emotions in controlled settings without much noise/interference, 2) focused on English speaking demographic and 3) hard labels for a speech segment [3]. In this work, we move away from the emotion recognition formulation and instead focus on modeling a more complex task of the ratio of people perceiving the given speech for every emotion. We utilize the first-ever speech-emotion corpus containing the fraction of the annotator population who categorize it as each of the nine emotions - Anger, Boredom, Calmness, Concentration, Determination, Excitement, Interest, Sadness, and Tiredness. This dataset is released as part of the Computational Paralinguistics Challenge (ComParE) 2023 [16]. It contains data from multi-lingual speakers in real-world settings making the task more challenging. We

want to design a robust architecture that can predict the 'share' of people who can perceive a given speech segment for each type of emotion. We provide a detailed analysis of the impact of different SSL embeddings on non-semantic downstream tasks like emotion perception. Non-semantic tasks rely on speech beyond its lexical meaning like language identification, speaker identification, emotion-related tasks, etc. We also discuss various architectures using custom attention, temporal convolution layers, long-short-term memory (LSTMs) and transformers, to learn a regressor for learning a continuous target for each emotion and provide insights.

The remainder of this paper is organized as follows. Section 2 describes the data preparation, the preliminaries on SSL models, and the implementation of various regressor architectures. Section 3 discusses the various evaluation settings that allow us to discern our findings and present them in a comprehensive manner. Finally, Section 4 summarises our contributions and discusses our future directions.

2 EXPERIMENTAL SET UP

We obtain the data from Hume AI as part of the Emotion Share Sub-Challenge of ACM Multimedia 2023 COMPUTATIONAL PARALINGUISTICS CHALLENGE (ComParE) [16]. More details on the data collection and statistics can be found in the baseline paper. This is a multi-lingual dataset with a wide demographic of participants. 9 different emotions are considered for each sample and are given a 'share' (based on the number of evaluators who rated the sample as a corresponding emotion) as the target [5]. This is a multi-task regression where we want to predict a continuous target for each sample for every emotion.

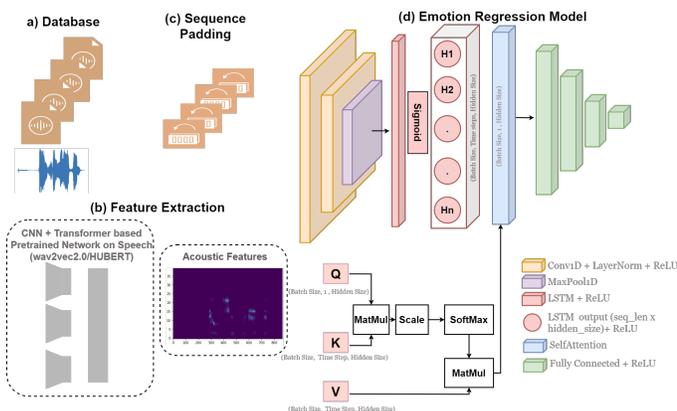


Figure 1: Overall architecture of our approach. a) Accessing the .wav files, b) Extracting embeddings from pretrained models or using acoustic features, c) Preparing dataset for training/evaluation by zero-padding the extracted features to deal with sequences of varying length, d) Using embedding/acoustic features for emotion share prediction.

2.1 Data Preparation

The audio files vary in size from 1s to 8s. To support different audio lengths we use the technique of padding and masking in the following sequence models. Instead of masking the raw audio-files to the length of the longest audio segment, we conduct feature-padding as shown in Fig. 1.(c). After we extract the features using

pretrained models on speech and acoustic features, we pad zeros to the temporal axis for each sample to match the maximum sequence length. In the case of wav2vec2.0 and HuBERT embeddings (both BASE and LARGE) we have a maximum sequence length of 398. For acoustic features, our maximum sequence length is 797. We store a dictionary of all the original sequence lengths as well, so that deeper in the architecture when operating on representations, we only consider the non-padded features.

2.2 Pretrained Self-Supervised Speech Features

Transformer-based large language models are the most sought after techniques to derive feature embeddings for a given speech segment. In this work, we explore two state-of-the-art self-supervised pipelines for feature extraction, wav2vec2.0 [2] and HuBERT [8]. The central architecture of both these models looks very similar superficially, consisting of convolution layers to extract latent feature embeddings from the raw audio followed by transformer layers. However, their training procedures are completely different. wav2vec2.0 uses contrastive loss by masking the latent feature embeddings and measuring the similarity between the predicted latent embedding versus the original feature embedding. They also leverage techniques like quantization and diversity loss to ensure robustness and avoid learning oversimplified representations. HuBERT on the other hand clusters the audio segments based on their latent feature embeddings using the K-Means algorithm and optimizes a cross-entropy loss. It chases the idea of discovering hidden units in a language rather than obtaining granular word-level representations as in the case of wav2vec2.0.

Both of these models are available in two sizes - BASE (90M parameter) and LARGE (300M parameters), with embedding sizes of 768 and 1024 respectively. Most previous works on emotion tasks are formulated as recognition. In this study, we address the unique task of predicting a continuous emotion perception target for a given audio segment for 9 different emotions. This is a more complex task than the former.

2.3 Acoustic Features

In many emotion classification tasks, a standard set of speech features such as Mel spectrogram, Mel-Frequency Cepstral Coefficients (MFCC), and raw spectrogram are studied [15]. Other non-semantic speech tasks like disfluency detection have also used MFCC and filter banks as their baseline method [13]. We have implemented acoustic feature extraction using a 40-dimensional mel-filterbank with cut-off frequencies at 0 Hz and 8000 Hz with a 25 ms window.

2.4 Emotion Regression Networks

After the embeddings have been extracted from the pre-trained models or using acoustic features, we use different types of regression networks to predict emotion share. In this subsection, we discuss different regression models. This part of the model is essential because the pre-trained self-supervised speech features are obtained from the models which are not specifically trained for emotion detection/share problems. For all the models discussed below, let us consider that the output from the pre-trained model (or acoustic features) is of dimension (N, L, W) where N , L , and W represents the batch size, sequence length, and feature size respectively.

Architecture 1: Architecture 1 comprises 1D CNN, an LSTM layer, and a feed-forward neural network (FFNN). 1D CNN layers increase

Pretrained model / Acoustic Features	Emotion Regression model	Anger	Boredom	Calmness	Concentration	Determination	Excitement	Interest	Sadness	Tiredness
wav2vec2.0 (Base)	Architecture 1	0.328	0.466	0.509	0.459	0.442	0.382	0.310	0.394	0.468
	Architecture 2	0.339	0.47	0.520	0.457	0.462	0.394	0.329	0.394	0.468
wav2vec2.0 (Large)	Architecture 1	0.138	0.129	0.266	0.333	0.311	0.278	0.156	0.230	0.407
	Architecture 2	0.125	0.261	0.353	0.392	0.206	0.235	0.141	0.097	0.413
HuBERT (Base)	Architecture 1	0.418	0.545	0.566	0.526	0.526	0.441	0.403	0.478	0.544
	Architecture 2	0.441	0.552	0.567	0.526	0.532	0.455	0.415	0.485	0.546
HuBERT (Large)	Architecture 1	0.450	0.566	0.583	0.538	0.539	0.485	0.427	0.505	0.560
	Architecture 2	0.467	0.575	0.587	0.542	0.547	0.476	0.426	0.518	0.557
MFCC Features	Architecture 1	0.231	0.358	0.404	0.290	0.363	0.301	0.148	0.125	0.390
	Architecture 2	0.180	0.343	0.379	-0.017	0.327	0.295	0.162	0.195	0.342
Baseline [16]		0.428	0.545	0.559	0.524	0.531	0.453	0.431	0.476	0.55

Table 1: Summary of Spearman correlation (ρ) for Architecture 1: CNN+LSTM+FFNN, and Architecture 2: CNN+LSTM+Attention+FFNN with wav2vec2.0 and HuBERT foundation models as feature extractors. HuBERT-Large with Architecture 2 provides the best performance.

the feature size from W to W_c but help in summarising the time series data and reduce the sequence length to $\frac{L}{2}$. The output of the 1D CNN layers is fed into a 2-layer LSTM network to capture any time dependencies in the signal in the context of the emotion share task. The output of the LSTM layer is of the dimension (N, L, W_l) . In the next step, we take the mean of the feature vector across all the time points of the LSTM output to obtain a tensor of dimension (N, W_l) . For a sample in a batch with L_i sequence length, the mean vector (\bar{f}_l) of size W_l can be obtained using Eq.1. In Eq.1, f_l^i represents the LSTM output vector at i^{th} time step. The mean output is then fed into FFNN to get the emotion share. It is important to note that the emotion regression network is different for all the emotions and they are trained separately too.

$$\bar{f}_l = \frac{\sum_{i=1}^{L_i} f_l^i}{L_i} \quad (1)$$

Architecture 2 : Architecture 2 is same as the Architecture 1 till the LSTM layer. In Architecture 2, after LSTM, we introduce the self-attention mechanism. Past works [17] demonstrate that self-attention improves performance in emotion detection tasks. We calculate scaled dot-product attention [20] using Query (Q), Key (K), and Value (V). In our model, we consider K and V to be the same as the LSTM output. Query vector (Q) is calculated from the final hidden state (h_n) of the LSTM. h_n has the dimension of $(2, N, W_l)$ since we are using 2 hidden layers. To obtain Q from the h_n , we first transform h_n to the dimension $(N, 2*W_l)$. Let us represent the transformed h_n as h'_n . Q is obtained from h'_n using a simple matrix multiplication shown in Eq.2a. In Eq.2a, the dimension of W_Q is $(2*W_l, W_l)$. In the next step, the attention (A) is calculated using Eq.2b. Attention A is then fed into FFNN to predict the emotion share. Fig. 1.(c) illustrates this architecture with pseudo-dimensions.

$$Q = \text{matmul}(W_Q, h'_n) \quad (2a)$$

$$A = \text{softmax}\left(\frac{QK^T}{W_l}\right)V \quad (2b)$$

2.5 Training Details

We run the experiments on an Ubuntu OS server equipped with NVIDIA TITAN RTX GPUs with PyTorch framework. We use the Adam optimizer to train the model with mini-batch gradient descent. We use early stopping based on validation data for the final model

selection. We have used 128 as the batch size with $1e-4$ as the rate of learning. More details on our hyperparameter settings can be found in the released codebase¹.

3 EXPERIMENTAL RESULTS

In this section, we discuss results obtained using embeddings from different pretrained models and different emotion regression networks discussed in the above sections. To assess the quality of fit we use Spearman correlation (ρ) [18] as it is used in the baseline paper [16] too. Spearman correlation indicates how well can the relationship between the two variables be described using a monotonic function and it is a favored metric for ranking-based values.

3.1 Evaluation on different Speech Embeddings

Table 1 shows the ρ value for all 9 emotions obtained from all the embeddings and acoustic features. Table 1 clearly shows that HuBERT-Large embeddings are the most effective. To understand the effect each embedding/acoustic feature has on the prediction, we show the percentage change in the average ρ (among 9 emotions) value with respect to the baseline model [16] in Table 2. It can be clearly noted from Table 1 that the HuBERT-Large embeddings outperform all other embeddings as well as acoustic features. Also, HuBERT-Large embeddings outperform the baseline by 4.6%. It is important to note that the **current task of predicting the emotion share is complex** as the ratings are dependent on human annotators. If we change the group of annotators, **the ground truth might vary**. Additionally, the data originates from **multi-lingual sources** but the pretrained models are trained only in the English language, making the task further complex. Therefore, for such a complex and uncertain task, **an improvement of 4.6%** is significant.

Although HuBERT-Large embedding performs the best, it does not indicate that using a large embedding size necessarily outperforms a smaller counterpart. This is evident from the fact that the wav2vec2.0-Base model outperforms the wav2vec2.0-Large model by 55%. **HuBERT (base and Large) embeddings outperform all other models as it is able to capture more non-semantic elements in the audio due to their training method.** They inherently capture more abstract representations of a language which is possibly the reason for their better performance over wav2vec2.0.

¹https://github.com/payalmohapatra/EmotionShare_ACM23.git

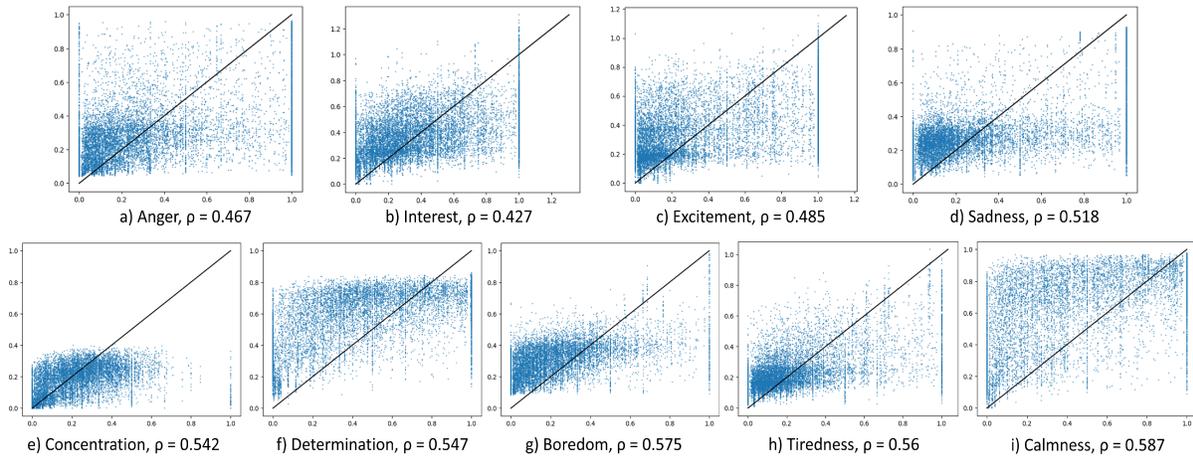


Figure 2: Scatter plot to show the quality of fit. The x-axis is the true value and the y-axis is the predicted value. The black line is a $y=x$ line; points lying on it are the perfect prediction.

Fig.2 shows the scatter plot for all emotions on the *dev* dataset. It reveals the regression imbalance in the data for most of the emotions. This can be one of the reasons for the correlation (ρ) to be low in our model as well as in the baseline model. The current performance can be further improved by adopting techniques to address the imbalance in regression data rather than feature extraction or updated architectures.

Pretrained model / Acoustic Features	Average ρ (a)	Baseline Average ρ (b)	Percentage Change (%) : $\frac{100*(a-b)}{b}$
wav2vec2.0 (Base)	0.425	0.500	-15
wav2vec2.0 (Large)	0.274		-45.2
HuBERT (Base)	0.502	0.500	0.4
HuBERT (Large)	0.523		4.6
Acoustic Features	0.299		-40.2

Table 2: Percentage change in average ρ for different embeddings/acoustic features compared to the baseline on *Dev* dataset. HuBERT-Large gives the most improvement.

3.2 Evaluation on different Emotion Regressor networks

In the overall model architecture, the emotion regressor network is the segment of the model which learns the paralinguistic part of the speech. Hence, it is extremely important to tune the network well. To that extent, as described in Sec.2.4, we are using two types of emotion regressor networks. The only difference between Architecture 1 and 2 is the presence of a self-attention mechanism in Architecture 2. Additionally, we have also done some experiments with a multi-headed transformer network [20], which performed very poorly so we exclude it from further analyses. Fig.3 shows the average ρ value for architecture 1 and 2 for different embeddings/acoustic features. It can be observed from Fig.3 that for wav2vec2.0-Base, HuBERT-base, and HuBERT-Large, Architecture 2 performs better by 1-1.8 %. While for wav2vec2.0-Large and acoustic features, Architecture 2 performs worse. These observations signify that **attention can improve the model** but if the initial embeddings/features are not good then attention just hurts the model by increasing the parameters.

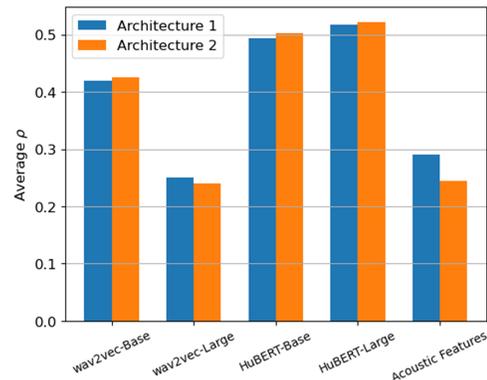


Figure 3: Comparing the results from Architecture 1 and 2.

4 CONCLUSION

In this work, we view speech emotion understanding as a perception by formulating a multi-label regression task rather than the commonly studied speech emotion recognition, a classification task. We demonstrate the effectiveness of our approaches on a unique dataset corpus with multilingual speakers to identify the share of people who perceive a given speech as one of the nine prominent emotions. We provide insights into using different foundation models as feature extractors and the role of attention in building an emotion regressor. Emotion relies on the non-semantic or paralinguistic aspects of speech. HuBERT-LARGE embeddings followed by a self-attention-based sequence model provide the best performance. In the Future, we are interested in addressing the current bottlenecks - 1) the use of foundation models pretrained only on English speech, 2) deep imbalance in the data for regression tasks, and 3) unavailability of meta-data like language identifiers, speaker characteristics(age, gender, etc.) as additional features.

REFERENCES

- [1] Hadhami Aouani and Yassine Ben Ayed. 2020. Speech emotion recognition with deep learning. *Procedia Computer Science* 176 (2020), 251–260.

- [2] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems* 33 (2020), 12449–12460.
- [3] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. *Language resources and evaluation* 42 (2008), 335–359.
- [4] Xingyu Cai, Jiahong Yuan, Renjie Zheng, Liang Huang, and Kenneth Church. 2021. Speech emotion recognition with multi-task learning. In *Interspeech*, Vol. 2021. 4508–4512.
- [5] Alan S Cowen, Hillary Anger Elfenbein, Petri Laukka, and Dacher Keltner. 2019. Mapping 24 emotions conveyed by brief human vocalization. *American Psychologist* 74, 6 (2019), 698.
- [6] Alan S Cowen and Dacher Keltner. 2021. Semantic space theory: A computational approach to emotion. *Trends in Cognitive Sciences* 25, 2 (2021), 124–136.
- [7] Pavol Harár, Radim Burget, and Malay Kishore Dutta. 2017. Speech emotion recognition with deep learning. In *2017 4th International conference on signal processing and integrated networks (SPIN)*. IEEE, 137–140.
- [8] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 29 (2021), 3451–3460.
- [9] Rashid Jahangir, Ying Wah Teh, Faiqa Hanif, and Ghulam Mujtaba. 2021. Deep learning approaches for speech emotion recognition: State of the art and research challenges. *Multimedia Tools and Applications* (2021), 1–68.
- [10] Sofoklis Kakouros, Themos Stafylakis, Ladislav Mošner, and Lukáš Burget. 2023. Speech-based emotion recognition with self-supervised models using attentive channel-wise correlations and label smoothing. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.
- [11] S Lalitha, D Geyasruti, R Narayanan, and M Shrivani. 2015. Emotion detection using MFCC and cepstrum features. *Procedia Computer Science* 70 (2015), 29–35.
- [12] MS Likitha, Sri Raksha R Gupta, K Hasitha, and A Upendra Raju. 2017. Speech based human emotion recognition using MFCC. In *2017 international conference on wireless communications, signal processing and networking (WiSPNET)*. IEEE, 2257–2260.
- [13] Payal Mohapatra, Bashima Islam, Md Tamzeed Islam, Ruochen Jiao, and Qi Zhu. 2023. Efficient Stuttering Event Detection Using Siamese Networks. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.
- [14] Edmilson Morais, Ron Hoory, Weizhong Zhu, Itai Gat, Matheus Damasceno, and Hagai Aronowitz. 2022. Speech emotion recognition using self-supervised features. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 6922–6926.
- [15] Sandeep Kumar Pandey, Hanumant Singh Shekhawat, and SR Mahadeva Prasanna. 2019. Deep learning techniques for speech emotion recognition: A review. In *2019 29th International Conference Radioelektronika (RADIOELEKTRONIKA)*. IEEE, 1–6.
- [16] Björn W Schuller, Anton Batliner, Shahin Amiriparian, Alexander Barnhill, Maurice Gerczuk, Andreas Triantafyllopoulos, Alice Baird, Panagiotis Tzirakis, Chris Gagne, Alan S Cowen, et al. 2023. The ACM Multimedia 2023 Computational Paralinguistics Challenge: Emotion Share & Requests. *arXiv preprint arXiv:2304.14882* (2023).
- [17] Jagjeet Singh, Lakshmi Babu Saheer, and Oliver Faust. 2023. Speech Emotion Recognition Using Attention Model. *International Journal of Environmental Research and Public Health* 20, 6 (2023). <https://doi.org/10.3390/ijerph20065140>
- [18] Charles Spearman. 1904. The Proof and Measurement of Association between Two Things. *The American Journal of Psychology* 15 [1], 72–101 (1904).
- [19] Andreas Triantafyllopoulos, Björn W Schuller, Gökçe İymen, Metin Sezgin, Xi-anheng He, Ziji Yang, Panagiotis Tzirakis, Shuo Liu, Silvan Mertes, Elisabeth André, et al. 2023. An overview of affective speech synthesis and conversion in the deep learning era. *Proc. IEEE* (2023).
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
- [21] Johannes Wagner, Andreas Triantafyllopoulos, Hagen Wierstorf, Maximilian Schmitt, Felix Burkhardt, Florian Eyben, and Björn W Schuller. 2023. Dawn of the transformer era in speech emotion recognition: closing the valence gap. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [22] Janghoon Yang. 2023. Ensemble deep learning with HuBERT for speech emotion recognition. In *2023 IEEE 17th International Conference on Semantic Computing (ICSC)*. IEEE, 153–154.