

Multimodal Approach for Assessing Neuromotor Coordination in Schizophrenia Using Convolutional Neural Networks

Yashish M. Siriwardena*
 Carol Espy-Wilson
 yashish@umd.edu
 espy@umd.edu
 University of Maryland
 College park, MD, USA

Chris Kitchen
 Deanna L. Kelly
 Ckitchen@jhu.edu
 dlkelly@som.umaryland.edu
 University of Maryland School of Medicine
 Baltimore, MD, USA

ABSTRACT

This study investigates the speech articulatory coordination in schizophrenia subjects exhibiting strong positive symptoms (e.g. hallucinations and delusions), using two distinct channel-delay correlation methods. We show that the schizophrenic subjects with strong positive symptoms and who are markedly ill pose complex articulatory coordination pattern in facial and speech gestures than what is observed in healthy subjects. This distinction in speech coordination pattern is used to train a multimodal convolutional neural network (CNN) which uses video and audio data during speech to distinguish schizophrenic patients with strong positive symptoms from healthy subjects. We also show that the vocal tract variables (TVs) which correspond to place of articulation and glottal source outperform the Mel-frequency Cepstral Coefficients (MFCCs) when fused with Facial Action Units (FAUs) in the proposed multimodal network. For the clinical dataset we collected, our best performing multimodal network improves the mean F1 score for detecting schizophrenia by around 18% with respect to the full vocal tract coordination (FVTC) baseline method implemented with fusing FAUs and MFCCs.

CCS CONCEPTS

• **Computing methodologies** → **Neural networks**; • **Social and professional topics** → **People with disabilities**.

KEYWORDS

Multimodal system, Schizophrenia, Articulatory coordination, FAUs, Vocal tract Variables

ACM Reference Format:

Yashish M. Siriwardena, Carol Espy-Wilson, Chris Kitchen, and Deanna L. Kelly. 2021. Multimodal Approach for Assessing Neuromotor Coordination in Schizophrenia Using Convolutional Neural Networks. In *Proceedings of the 2021 International Conference on Multimodal Interaction (ICMI '21)*, October 18–22, 2021, Montréal, QC, Canada. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3462244.3479967>

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 ACM 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](https://permissions.acm.org).

ICMI '21, October 18–22, 2021, Montréal, QC, Canada

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8481-0/21/10...\$15.00

<https://doi.org/10.1145/3462244.3479967>

1 INTRODUCTION

Schizophrenia is a chronic mental disorder with heterogeneous presentations that affect around 60 million (1%) of the world's adult population [12]. Symptoms of schizophrenia are broadly categorized as either positive, which are pathological functions not present in healthy individuals (e.g., hallucinations and delusions); negative, which involve the loss of functions or abilities (e.g., apathy, lack of pleasure, blunted affect and poor thinking); or cognitive (deficits in attention, memory and executive functioning) [2, 4]. Previous studies have shown promising results in identifying the severity of depression by using coordination features derived from the correlation structure of the movements of various articulators [6]. Based on this, a preliminary study was done by Siriwardena et al. [18] to understand how positive symptoms of schizophrenia affect the articulatory coordination in speech. These findings are the impetus for the current investigation where more subjects and data are used to validate the fact that neuromotor coordination is altered in schizophrenic patients who are markedly ill and exhibit strong positive symptoms.

Time-delay embedded correlation (TDEC) analysis has shown promising results in assessing neuromotor coordination in Major Depressive Disorder (MDD), and the eigenspectra derived from the correlation matrices have been used effectively for classification of MDD subjects from healthy [17, 22, 24]. Recently, new multi-scale full vocal tract coordination (FVTC) features generated with a dilated CNN have shown further improvement in classification for selected datasets of MDD subjects [9]. The FVTC method addresses repetitive sampling and matrix discontinuity issues of TDEC analysis by introducing a new channel-delay correlation matrix. In this paper, we compare both TDEC and FVTC methods for generating input correlation matrices for training a multimodal CNN with audio and video features. We also propose a model which uses both TDEC and FVTC correlation structures to classify subjects with strong positive symptoms in schizophrenia from healthy.

Throughout the paper, we compare and contrast studies mostly done in depression analysis since it is hard to find any comparable speech articulatory coordination based analysis on schizophrenia. While there are studies investigating the changes in language used when subjects are schizophrenic with positive symptoms [13], to our knowledge this study is one of the first papers to present results that show changes in the coordination of speech gestures produced by schizophrenic subjects. Extending the results in [18], this paper also validates that schizophrenic subjects with strong positive symptoms have a more complex articulatory coordination with

Table 1: Details of the dataset used

	SZ	HC
Number of subjects	7	11
Number of sessions	17	34
Mean session duration	35min	18min
Number of utterances	1208	1132
Hours of speech	10.0	9.43

respect to healthy controls. Finally, this study presents the importance and effectiveness of multimodal fusion (audio and video) for screening mental health disorders like schizophrenia by proposing a new CNN based multimodal architecture.

2 DATABASE AND FEATURES

2.1 Database

A database recently collected for a collaborative observational study conducted by the University of Maryland School of Medicine and the University of Maryland College Park has been used for this study [11]. The database contains video and audio data of free response assessments administered in an interview format. Data for this study was collected from 23 schizophrenic (SZ) patients and 20 healthy controls (HC). All of the schizophrenic patients were clinically diagnosed. Every subject participated in four interview sessions over a period of six weeks. Each interview session is 10-45 minutes long and every subject is assessed using standard depression severity measures and global psychopathology measures by a clinician and themselves. For this study, we used the clinician assessments based on the 18-item Brief Psychiatric Rating Scale (BPRS) [10], where we selected subjects based on the total BPRS score, and the subscores for psychosis (BPRS item11, item12, item4, item15) and activation (BPRS item6, item7, item17), and the Hamilton Rating Scale for Depression (HAMD) [8].

Table 1 presents the information on the dataset used for the study. The 7 schizophrenic subjects (4 Males, 3 Females) are selected such that they are markedly ill (BPRS total ≥ 45), have higher sub scores for psychosis and activation, but are not depressed or only mildly depressed (HAMD between 0 and 14). The 11 healthy controls (5 Males, 6 Females) are chosen such that they are not depressed (HAMD < 7) or schizophrenic (BPRS < 32).

The audio data were first diarized using transcripts which include the speaker ID and time stamps to separate out the utterances which correspond to the subject from the interviewer. The utterances which are longer than 40 seconds were then segmented into 40 second chunks. If the remaining amount was less than 5 seconds (the minimum length accepted), then it was added back to the last segment. Thus, for all the classification experiments, we used segments with a minimum length of 5 seconds and a maximum length of 45 seconds.

2.2 Facial Action Units (FAUs)

The video-based Facial Action Units (FAUs) provide a formalized method for identifying changes in facial expressions. We used the

Openface 2.0: Facial Behaviour Analysis toolkit [3] to extract seventeen FAUs from the recorded videos of the subjects during the interviews. The FAU features were sampled at a rate of 28 frames per second. We only analyzed those portions of the video when the subject was talking. For parallel delay CNN model in section 3.1 and FVTC model in section 3.2, we choose only 10 FAUs (FAU 6,7,9,10,12, 14,15,17,20 and 23 from FACS coding system [14]) which are near the mouth area that can capture coordination of lip and near lip movements during the voice activity. One other reason to choose 10 FAUs is to handle the high dimensionality of the TDEC correlations structure which poses computational limitation when training multi-modal networks in section 4.3.

2.3 Vocal Tract Variables (TVs)

We used a speech inversion (SI) system [19, 20] that maps the acoustic signal into 6 vocal tract variables (TVs). The 6 TVs are namely Lip Aperture (LA), Lip Protrusion (LP), Tongue Body Constriction Location (TBCL), Tongue Body Constriction Degree (TBCD), Tongue Tip Constriction Location (TTCL) and, Tongue Tip Constriction Degree (TTCD).

Seneviratne et al. [17] in a recent study showed that incorporating glottal TVs generated by periodicity and aperiodicity measures (by Aperiodicity, Periodicity and Pitch (APP) detector [5]) improved the results of depression detection. Thus, in this study we use 6 TVs generated from the SI along with the 2 glottal TVs as the key audio features for the classification models.

2.4 Mel-Frequency Cepstral Coefficients (MFCCs)

Previous studies in depression prediction using speech [15, 16] have shown the superiority of MFCCs over other audio based features like extended Geneva Minimalistic Acoustic Parameter Set (eGeMAPS) [7] and DEEP SPECTRUM features [1]. Huang et al. [9] showed with their depression classification study that coordination features computed from MFCCs perform better with respect to formants and eGeMAPS features. So to compare how robust and effective the TVs are for detecting schizophrenia, we chose MFCCs as the baseline audio features for our study. We extracted 13 MFCCs from the librosa python library using an analysis window of 20 ms with a 10 ms frame shift. Only 12 MFCCs were used for analysis by discarding the 1st coefficient.

3 METHODOLOGY

3.1 Time-delay embedded correlation Analysis (TDEC)

Coordination among 10 FAUs, 6 TVs and the 12 MFCCs were estimated using the correlation structure features. The features are estimated by computing a channel delay correlation matrix using time delay embedding at two delay scales [6, 23, 24]. The computed correlation matrix can be considered as a compact representation which provides rich information on the underlying mechanisms in articulatory coordination. Each correlation matrix R_i has a dimensionality (MN x MN) where M = 10, 8 and 12 for FAUs, TVs and MFCCs respectively. N corresponds to the number of time delays per channel which is 15 for all the considered feature types.

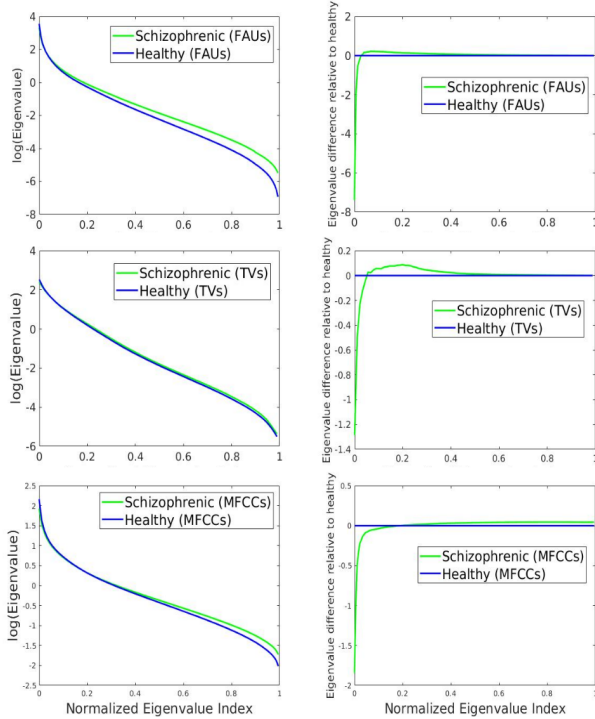


Figure 1: Averaged eigenspectra (on left) and the corresponding difference plots (on right) for FAUs, MFCCs and TVs

From the correlation matrix R_i calculated for each sample i , the eigenspectrum is computed. The eigenspectrum generated for FAUs is a 150-dimensional vector which is rank ordered (in descending order of magnitude of eigenvalues) from index $j=1, \dots, 150$. The eigenspectra generated from TVs and MFCCs are 120-dimensional and 180-dimensional, respectively.

Figure 1 shows the averaged eigenspectra (on left) computed for FAUs, TVs and MFCCs. The difference plots in Figure 1 (in right) are calculated by taking the difference between averaged eigenspectra curve for speech from schizophrenic subjects with respect to that of healthy controls. These eigenspectra and difference plots help us to understand the coordination complexity of the speech and facial gestures. The magnitude of each eigenvalue is proportional to the amount of correlation in the direction of their associated eigenvectors. The difference plots show that schizophrenic speech has smaller low-rank eigenvalues relative to the healthy controls, and the trend is reversed towards the high-rank eigenvalues. Therefore, schizophrenic speech needs a larger number of independent dimensions implying a more complex articulatory coordination than speech from healthy controls [21, 23]. The same argument is true for facial gestures.

3.2 Full vocal tract coordination (FVTC)

Huang et al. [9] in a recent study with MDD introduces a new channel delay correlation method inspired by TDEC, which uses a different correlation structure with correlations starting from 0 to a

delay of 'D' frames (a design choice). The delayed autocorrelations and cross-correlations across channels are stacked to form the FVTC correlation structure. FVTC includes every correlation within the considered D frames and also avoids the repetitive use of same correlations as in the TDEC correlation matrix.

4 MULTIMODAL SYSTEMS

4.1 Parallel delay scale TDEC-CNN model (TDEC-CNN) : Model 1

We developed a CNN architecture which takes in multiple time-delay embedded correlation matrices with 2 delay scales in parallel as inputs for two 2D-CNN layers. The output from the 2 CNN layers are then concatenated and passed through another 2D-CNN layer. Batch normalization, max-pooling and dropout were applied afterwards. The flattened output is then fed to a fully connected layer with 64 neurons. 16 filters with kernel size (3,3) was used for every 2D CNN layer and every CNN layer has ReLU activation. We trained individual models for FAUs, TVs and MFCCs where 3 and 7 sample delay scales were used for FAUs and 7 and 15 sample delay scales were used for TVs and MFCCs.

4.2 FVTC CNN model (FVTC-CNN) : Model 2

We designed a CNN model inspired by the one in [9] which takes the FVTC correlation matrix computed in section 3.2 as the input. To reduce the number of trainable parameters in the original model [9], we reduced the size of the two fully connected layers to 64 and 8. We used the same dilation rates 1,3,7,15 as in the original model. We chose 45 as the 'D' parameter for FAUs and 50 for TVs and MFCCs. The 'D' values were chosen from the set of (45,50,55) after doing a grid search on individual feature based systems.

4.3 Multimodal CNNs

One of the key contributions of the paper comes with the multimodal networks developed to fuse the audio and video based coordination features calculated from the TDEC and FVTC methods. 3 types of fusion models were designed namely the (Parallel delay TDEC-FVTC CNN), (FVTC-FVTC CNN) and (Parallel delay TDEC- Parallel delay TDEC CNN). The same model architectures developed in section 4.1 and section 4.2 were used and they are fused after passing through all the 2D-CNN layers by concatenating the flattened outputs from each model. The concatenated output then passes through 2 fully connected layers which have 64 and 8 neurons, respectively. That output is then fed to a softmax output layer to generate the final class probabilities (for schizophrenic and healthy classes).

To choose the learning rate and batch size for the individual and multimodal systems, we did a grid search from sets of (1e-4, 3e-4, 1e-5, 3e-5) and (32,64,128) for learning rates and batch sizes, respectively. 1e-4 for learning rate and 64 for batch size gave the best metrics in classification for the best performing model in Table 3. Every model was trained in leave one subject out cross validation fashion where both accuracy and F1 scores are calculated across 18 folds. This ensured that every model is always trained with around 2000 sample points (17 subject's data) and then tested on the excluded subject (test fold). To come up with the subject level prediction label from the multiple segment level predictions, we

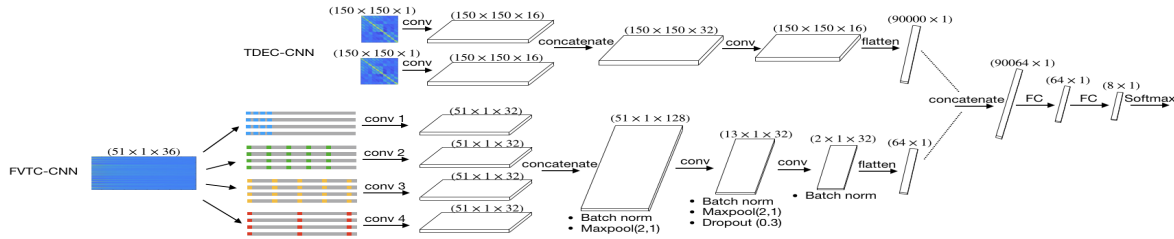


Figure 2: TDEC and FVTC combined multimodal architecture for best performing model in Table 3

Table 2: Individual Model results for FAUs, TVs and MFCCs. Best Model for each feature type is highlighted in bold.

Features	TDEC-CNN (Model1)		FVTC-CNN (Model2)	
	Accuracy	F1(S)/F1(H)	Accuracy	F1(S)/F1(H)
FAU	83.33%	0.80/0.86	83.33%	0.77/0.87
TV	66.67%	0.57/0.73	72.22%	0.62/0.78
MFCC	61.11%	0.46/0.70	72.22%	0.55/0.80
MFCC+Glottal TVs	60.05 %	0.45/0.69	72.22%	0.55/0.80

use the best 25% of the total segments in the test fold which are predicted with the highest class probabilities (even if the model predicts the wrong class). Every model is trained for 200 epochs with early stopping based on validation loss with patience of 15 epochs. Figure 2 shows the architecture for the best performing multimodal system from Table 3.

Table 3: Multimodal classification results

Models	Accuracy	F1(S)/F1(H)
FAU (Model2)+TV(Model2)	66.67%	0.67/0.67
FAU (Model1)+TV(Model1)	72.22%	0.67/0.76
FAU (Model2)+MFCC(Model2)	72.22%	0.62/0.78
FAU (Model1)+MFCC(Model2)	77.78%	0.67/0.83
FAU (Model1)+(MFCC+Glottal TVs)(Model2)	83.33%	0.73/0.88
FAU (Model1)+TV(Model2)	88.89%	0.86/0.91

5 DISCUSSION

The averaged eigenspectra and difference plots in figure 1 shows that the low-rank eigenvalues are smaller for schizophrenic subjects relative to the healthy controls, and this trend is reversed towards the high-rank eigenvalues. A key observation associated with depression severity [6, 22, 24] is that low-rank eigenvalues are larger for MDD subjects relative to healthy controls where as they are smaller for high-rank eigenvalues. The magnitude of high-rank eigenvalues indicates the dimensionality of the time-delay embedded feature space. Thus, larger values in the high-rank eigenvalues can be associated with greater complexity of articulatory coordination[6, 23]. Therefore we can infer that the schizophrenic subjects with strong positive symptoms have a higher complexity than the healthy controls and the MDD patients. These results are likely due to the negative symptoms of depression which results

in psychomotor slowing (i.e., simpler coordination) and the strong positive symptoms of the schizophrenic patients such as activation that results in motor hyperactivity (i.e., complex coordination). Supporting our hypothesis, we see this effect from eigenvalues computed from the FAUs, TVs and MFCCs.

Table 2 shows the average accuracy across the 18 folds and the F1 scores for classifying schizophrenic and healthy subjects by training individual models for FAUs, TVs and MFCCs based on TDEC-CNN and FVTC-CNN models. Results suggest that FAUs perform the best when compared to TVs and MFCCs in classification metrics. This could be due to the inclusion of a wider range of facial muscle movements which are not limited to only those around the speech articulators. Moreover, FVTC-CNN models trained with TVs and MFCCs perform the best when compared to TDEC-CNN models trained with the same features. With respect to F1 scores, TVs perform better than the MFCCs in both FVTC-CNN and TDEC-CNN models showing the robustness of TVs in capturing the articulatory changes in speech.

Table 3 shows results for the 6 multimodal systems trained by fusing video and audio features from TDEC and FVTC methods. It is interesting to note that the models with heterogeneous architectures perform the best when compared to models which use the same correlation structure for both audio and video features. To do a fair comparison between TVs and MFCCs, we also trained individual and multimodal networks by incorporating the 2 glottal TVs along with the 12 MFCCs, so that glottal source level information is also accounted. Even then, Table 3 shows that the TV based best performing model out performs the MFCC based multimodal systems. Furthermore, the multimodal system which uses TDEC and FVTC correlation structures for FAUs and TVs, respectively, outperforms the baseline model trained on FVTC for both FAUs and MFCCs by around 18% in terms of the mean F1 score. This result is interesting in the sense that the video features complement well with TVs over MFCCs in the proposed multimodal setting.

In conclusion, this paper proposes a multimodal approach to classify subjects with strong positive symptoms in schizophrenia from healthy. We also show that the video based features are more effective in identifying articulatory coordination changes while also asserting the fact that fusing with audio based TVs significantly boost the performance in classification.

ACKNOWLEDGMENTS

This work was supported by a UMCP UMB - AI + Medicine for High Impact (AIM-HI) Challenge Award.

REFERENCES

- [1] Shahin Amiriparian, Maurice Gerczuk, Sandra Ottl, Nicholas Cummins, M. Freitag, Sergey Pugachevskiy, Alice Baird, and B. Schuller. 2017. Snore Sound Classification Using Image-Based Deep Spectrum Features. In *INTERSPEECH*.
- [2] Nancy C. Andreasen and Scott Olsen. 1982. Negative v Positive Schizophrenia: Definition and Validation. *Archives of General Psychiatry* 39, 7 (07 1982), 789–794. <https://doi.org/10.1001/archpsyc.1982.04290070025006>
- [3] T. Baltrusaitis, A. Zadeh, Y. C. Lim, and L. Morency. 2018. OpenFace 2.0: Facial Behavior Analysis Toolkit. In *2018 13th IEEE International Conference on Automatic Face Gesture Recognition (FG 2018)*. 59–66.
- [4] Caroline Demily and Nicolas Franck. 2008. Cognitive remediation: a promising tool for the treatment of schizophrenia. *Expert Review of Neurotherapeutics* 8, 7 (2008), 1029–1036. <https://doi.org/10.1586/14737175.8.7.1029> arXiv:<https://doi.org/10.1586/14737175.8.7.1029> PMID: 18590474.
- [5] O. Deshmukh, C. Y. Espy-Wilson, A. Salomon, and J. Singh. 2005. Use of temporal information: detection of periodicity, aperiodicity, and pitch in speech. *IEEE Transactions on Speech and Audio Processing* 13, 5 (2005), 776–786. <https://doi.org/10.1109/TSA.2005.851910>
- [6] Carol Espy-Wilson, Adam C. Lammert, Nadee Seneviratne, and Thomas F. Quatieri. 2019. Assessing Neuromotor Coordination in Depression Using Inverted Vocal Tract Variables. In *Proc. Interspeech 2019*. 1448–1452. <https://doi.org/10.21437/Interspeech.2019-1815>
- [7] F. Eyben, K. R. Scherer, B. W. Schuller, J. Sundberg, E. André, C. Busso, L. Y. Devillers, J. Epps, P. Laukka, S. S. Narayanan, and K. P. Truong. 2016. The Geneva Minimalistic Acoustic Parameter Set (GeMAPS) for Voice Research and Affective Computing. *IEEE Transactions on Affective Computing* 7, 2 (2016), 190–202. <https://doi.org/10.1109/TAFFC.2015.2457417>
- [8] Jeffrey S. Gonzalez, Erica Shreck, and Abigail Batchelder. 2013. *Hamilton Rating Scale for Depression (HAM-D)*. Springer New York, New York, NY, 887–888. https://doi.org/10.1007/978-1-4419-1005-9_198
- [9] Zhaocheng Huang, J. Epps, and D. Joachim. 2020. Exploiting Vocal Tract Coordination Using Dilated CNNs For Depression Detection In Naturalistic Environments. *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (2020)*, 6549–6553.
- [10] Edward E. Hunter and Meghan Murphy. 2011. *Brief Psychiatric Rating Scale*. Springer New York, New York, NY, 447–449. https://doi.org/10.1007/978-0-387-79948-3_1976
- [11] Deanna L. Kelly, Max Spaderna, Vedrana Hodzic, Suraj Nair, Christopher Kitchen, Anne E. Werkheiser, Megan M. Powell, Fang Liu, Glen Coppersmith, Shuo Chen, and Philip Resnik. 2020. Blinded Clinical Ratings of Social Media Data are Correlated with In-Person Clinical Ratings in Participants Diagnosed with Either Depression, Schizophrenia, or Healthy Controls. *Psychiatry Research* 294 (2020), 113496. <https://doi.org/10.1016/j.psychres.2020.113496>
- [12] G. R. Kuperberg. 2010. Language in schizophrenia Part 1: an Introduction. *Lang Linguist Compass* 4, 8 (Aug 2010), 576–589.
- [13] Fereshteh Momeni and Shahla Raghibdoust. 2012. The relationship between incoherent speech and different types of delusions and hallucinations in schizophrenics with positive symptoms. *Procedia - Social and Behavioral Sciences* 32 (2012), 288–295. <https://doi.org/10.1016/j.sbspro.2012.01.042> The 4th International Conference of Cognitive Science.
- [14] Emily B. Prince, Katherine B. Martin, and D. Messinger. 2015. Facial Action Coding System.
- [15] Anupama Ray, Siddharth Kumar, Rutvik Reddy, Prerana Mukherjee, and Ritu Garg. 2019. Multi-level Attention network using text, audio and video for Depression Prediction. arXiv:1909.01417 [cs.CV]
- [16] Fabien Ringeval, Björn Schuller, Michel Valstar, Nicholas Cummins, Roddy Cowie, Leili Tavabi, Maximilian Schmitt, Sina Alisamir, Shahin Amiriparian, Eva-Maria Messner, Siyang Song, Shuo Liu, Ziping Zhao, Adria Mallol-Ragolta, Zhao Ren, Mohammad Soleymani, and Maja Pantic. 2019. AVEC 2019 Workshop and Challenge: State-of-Mind, Detecting Depression with AI, and Cross-Cultural Affect Recognition. arXiv:1907.11510 [cs.HC]
- [17] Nadee Seneviratne, James R. Williamson, Adam C. Lammert, Thomas F. Quatieri, and Carol Espy-Wilson. 2020. Extended Study on the Use of Vocal Tract Variables to Quantify Neuromotor Coordination in Depression. In *Proc. Interspeech 2020*. 4551–4555. <https://doi.org/10.21437/Interspeech.2020-2758>
- [18] Yashish M. Siriwardena, Chris Kitchen, Deanna L. Kelly, and Carol Espy-Wilson. 2021. Inverted Vocal Tract Variables and Facial Action Units to Quantify Neuromotor Coordination in Schizophrenia. In *Proc. 12th International Seminar on Speech Production (ISSP 2020)*. 174–177. <https://issp2020.yale.edu/ProcISSP2020.pdf>
- [19] Ganesh Sivaraman. 2017. *Articulatory representations to address acoustic variability in speech*. Ph.D. Dissertation. University of Maryland College Park. <https://drum.lib.umd.edu/handle/1903/20422>
- [20] Ganesh Sivaraman, Vikramjit Mitra, Hosung Nam, Mark K. Tiede, and Carol Y. Espy-Wilson. 2016. Vocal Tract Length Normalization for Speaker Independent Acoustic-to-Articulatory Speech Inversion. In *Proceedings of Interspeech*. 455–459. <https://doi.org/10.21437/Interspeech.2016-1399>
- [21] James R. Williamson, Daniel W. Bliss, David W. Browne, and Jaishree T. Narayanan. 2012. Seizure prediction using EEG spatiotemporal correlation structure. *Epilepsy & Behavior* 25, 2 (2012), 230 – 238. <https://doi.org/10.1016/j.yebeh.2012.07.007>
- [22] James R. Williamson, Thomas F. Quatieri, Brian S. Helfer, Gregory Ciccarelli, and Daryush D. Mehta. 2014. Vocal and Facial Biomarkers of Depression Based on Motor Incoordination and Timing. In *Proceedings of the 4th International Workshop on Audio/Visual Emotion Challenge (Orlando, Florida, USA) (AVEC '14)*. Association for Computing Machinery, New York, NY, USA, 65–72. <https://doi.org/10.1145/2661806.2661809>
- [23] James R. Williamson, Thomas F. Quatieri, Brian S. Helfer, Rachele Horwitz, Bea Yu, and Daryush D. Mehta. 2013. Vocal Biomarkers of Depression Based on Motor Incoordination. In *Proceedings of the 3rd ACM International Workshop on Audio/Visual Emotion Challenge (Barcelona, Spain) (AVEC '13)*. Association for Computing Machinery, New York, NY, USA, 41–48. <https://doi.org/10.1145/2512530.2512531>
- [24] James R. Williamson, Diana Young, Andrew A. Nierenberg, James Niemi, Brian S. Helfer, and Thomas F. Quatieri. 2019. Tracking depression severity from audio and video based on speech articulatory coordination. *Computer Speech & Language* 55 (2019), 40 – 56. <https://doi.org/10.1016/j.csl.2018.08.004>