

# A CONTEXT-AWARE COMPUTATIONAL APPROACH FOR MEASURING VOCAL ENTRAINMENT IN DYADIC CONVERSATIONS

Rimita Lahiri<sup>1</sup>, Md Nasir<sup>2</sup>, Catherine Lord<sup>3</sup>, So Hyun Kim<sup>4</sup>, Shrikanth Narayanan<sup>1</sup>

<sup>1</sup>Signal Analysis and Interpretation Laboratory, University of Southern California, Los Angeles, USA

<sup>2</sup>Microsoft AI for Good Research Lab, Redmond, Washington, USA

<sup>3</sup>Semel Institute of Neuroscience and Human Behavior, University of California, Los Angeles, USA

<sup>4</sup>School of Psychology, Korea University, Seoul, Korea

## ABSTRACT

Vocal entrainment is a social adaptation mechanism in human interaction, knowledge of which can offer useful insights to an individual’s cognitive-behavioral characteristics. We propose a context-aware approach for measuring vocal entrainment in dyadic conversations. We use conformers (a combination of convolutional network and transformer) for capturing both short-term and long-term conversational context to model entrainment patterns in interactions across different domains. Specifically we use cross-subject attention layers to learn intra- as well as interpersonal signals from dyadic conversations. We first validate the proposed method based on classification experiments to distinguish between *real* (consistent) and *fake* (inconsistent/shuffled) conversations. Experimental results on interactions involving individuals with Autism Spectrum Disorder also show evidence of a statistically-significant association between the introduced entrainment measure and clinical scores relevant to symptoms, including across gender and age groups.

**Index Terms**— entrainment, context, transformers, convolution

## 1. INTRODUCTION

Interpersonal human interactions, notably dyadic interactions (interactions involving two people), are widely studied by social science and human-centered computing researchers alike [1, 2]. Such interactions are characterized by rich information exchange across multiple modalities including speech, language, and visual cues. Over the years, a significant amount of effort has been invested in developing tools for both conversational data collection and in understanding and modeling the signals extracted from these interactions.

A phenomenon called *entrainment* [3, 4] has been described as one of the major driving forces of an interaction [5]. While entrainment can be exhibited within and across different modalities, vocal entrainment [6] or acoustic-prosodic entrainment [7, 4, 8] is defined as an interlocutor’s tendency to accommodate or adapt to the vocal patterns of the other interlocutor over the course of the interaction. Understand-

ing entrainment [9] can provide meaningful insights to analyze behavioral characteristics of the individual interlocutors and the interaction participants. For example, a higher degree of entrainment is associated with positive behavioral markers like social desirability, smoother interactions, higher rapport content *etc.*[10, 11]. Entrainment can also serve as a valuable instrument to characterize behaviors in the study and practice of psychiatry and developmental studies involving distressed couples, children with autism spectrum disorder, addiction, *etc* [6, 9].

Due to the complex nature of entrainment and a scarcity of appropriately labeled speech corpora, quantifying entrainment is a challenging task. Most of the early works have relied on empirical and knowledge-driven tools like correlation, recurrence analysis, time-series analysis, spectral methods to measure how much a speaker is entraining to the other speaker [12]. This body of work often relied on the assumption of a linear relationship between the extracted entrainment representations and vocal features, which may not always hold. On the other hand, although context during a conversation plays an important part in interpersonal interactions, it has not been incorporated in existing approaches for measuring entrainment. While the recent line of works [6] employ a more direct data-driven strategy to extract entrainment related information from raw speech features, such are formulated in a way that they inherently only consider short-term context while overlooking more long-term context. Recently context-aware deep learning architectures such as transformers [13] have been proposed to capture richer contexts by explicitly modeling the temporal dimension and found many applications in natural language processing, speech and vision. In light of their success in modeling rich temporal context, we investigate if transformers can help capture meaningful information for quantifying entrainment.

In this work, we develop a context-aware model for computing entrainment, addressing the need for both short and long-range temporal context modeling. For the scope of this work, the proposed framework incorporates ‘context’ by aiming to train the model to learn the influence of the speakers

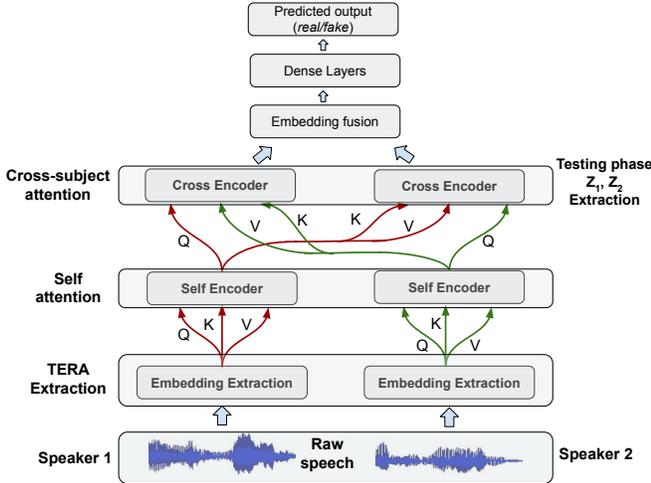


Fig. 1: Architecture for CED extraction.

on each other. We follow the established strategy of using a distance-based measure between consecutive turn-pairs in the projected embedding space and introduce the *Contextual Entrainment Distance (CED)* metric. The main contributions of this work are two fold: first, we use a combination of self-attention and convolution to extract both short-term and long-term contextual information related to entrainment; and second, we propose a transformer-based cross-subject framework for joint modeling of the interacting speakers to learn the pattern of entrainment. We experimentally evaluate the validity and efficacy of CED in dyadic conversations involving children and study its association with respect to different clinically-significant behavioral ratings where the role of entrainment has been previously implicated [9].

## 2. COMPUTING CONTEXT-AWARE ENTRAINMENT MEASURE

### 2.1. Unsupervised model training and CED computation

Prior literature in this domain have relied on computing a distance measure directly between the turn-level speech features  $X_1$  and  $X_2$  from speaker 1 and speaker 2 respectively [4]. However, these features also capture additional information such as speaker characteristics and ambient acoustic information which do not contribute towards learning the target entrainment patterns. The objective is to learn the inverse mapping between the embedding space  $(Z_1, Z_2)$  and the feature space  $(X_1, X_2)$  such that the model should learn to recognize turn pairs with high and low levels of entrainment.

Here, we formulate the problem by training the network to classify between interactions having consecutive turn segments (*true samples*) and interactions having random/shuffled turn segments (*fake samples*). We temporally partition the conversational audio sequence into speaker specific chunks and feed these chunks to the model to predict whether the fed audio chunks are part of real conversation or a fake one.

After the training phase, we use the trained network weights to extract the cross-encoder layer outputs for both speakers. Next, we calculate CED as the smooth L1 distance [6] between the embeddings obtained in the previous step.

### 2.2. Model architecture

As shown in Fig. 1, we use two main modules to build the model to compute entrainment, first, the self-attention encoder that is used to enhance the extracted features by attending to themselves and then, a cross-attention encoder which allows the features to attend to a different source.

We use conformer [14] layers for the self-attention module to model both short-term and long-term dependencies within an audio sequence in a parameter-efficient way by incorporating a convolutional module in the transformer layer. The self-attention layer obtains meaningful representation from the long-term interaction and the convolution layer is used to learn the local relation amongst the interaction based features.

To extract meaningful information related to entrainment, previous works have mostly relied on individual modeling of interlocutors involved in a conversation. However, entrainment being an interpersonal phenomenon, the need for jointly modeling interlocutors becomes heightened in such scenarios. We address this issue by using a transformer layer for cross-subject attention, allowing the features extracted per subject to access each other to capture crossed influence over the interaction.

## 3. EXPERIMENTS

### 3.1. Datasets

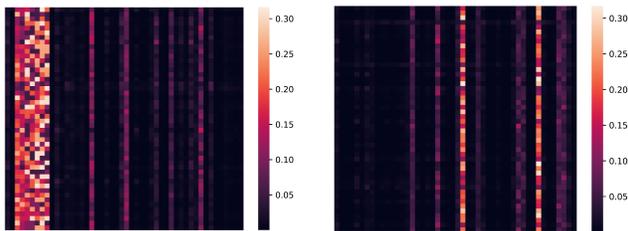
We use the following two datasets for our experiments.

*The Fisher Corpus* English Part 1 (LDC2004S13) [15] consists of spontaneous telephonic conversations between two native English speaking subjects. There are 5850 conversations of approximately 10 minutes duration. The dataset is accompanied with transcripts along with timestamps marking speaker duration boundaries. We use 60% of this dataset for training and 5% for testing.

*The ADOSmod3 corpus* consists of recorded conversations from autism diagnostic sessions between a child and a clinician who is trained to observe the behavioral traits of the child related to *Autism Spectrum Disorder (ASD)*. A typical interactive session following the *Autism Diagnostic Observation Schedule (ADOS)-2* instrument lasts about 40-60 minutes, and these sessions are composed of a variety of subtasks to evoke spontaneous response from the children under different social and communicative circumstances. In this work, we consider the administration of Module 3 meant for verbally fluent children and adolescents. Moreover, we focus on *Emotions* and *Social difficulties and annoyance* subtasks as these are expected to extract significant spontaneous speech and reaction from the child while answering questions

**Table 1:** Demographic details of ADOSMod3 dataset

Category	Statistics
Age(years)	Range: 3.58-13.17 (mean,std):(8.61,2.49)
Gender	123 male, 42 female
Non-verbal IQ	Range: 47-141 (mean,std):(96.01,18.79)
Clinical Diagnosis	86 ASD,42 ADHD
	14 mood/anxiety disorder
	12 language disorder
	10 intellectual disability, 1 no diagnosis
Age distribution	Cincinnati: $\leq 5$ yrs 7, 5-10 yrs 52, $\geq 10$ yrs 25 Michigan: $\leq 5$ yrs 11, 5-10 yrs 42, $\geq 10$ yrs 28



(a) Cross-encoder 1

(b) Cross-encoder 2

**Fig. 2:** Attention activations

**Table 2:** Classification experiment for *real vs fake* sessions

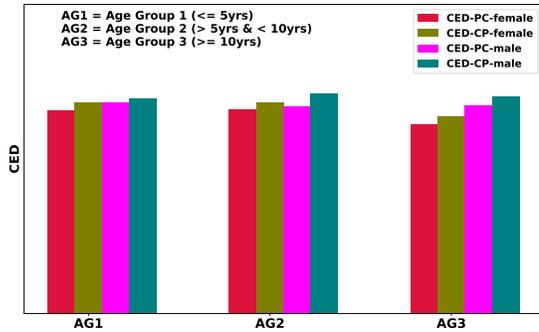
Measure	Classification accuracy(%)	
	Fisher Corpus	ADOSMod3 Corpus
Baseline 1	80.52	82.22
Baseline 2	76.33	70.64
Baseline 3	82.91	85.73
CED	92.13	95.66

about different emotions and social difficulties. The corpus consists of recordings from 165 children collected across 2 different clinical sites. We use this corpus for evaluation purpose, the demographic details of the dataset are reported in Table 1.

### 3.2. Experimental setup

#### 3.2.1. Feature extraction

In this work, to compute CED the speech segments of interest are conversational turns from both speakers. We compute the speaker turn boundaries from the time information available in the transcripts, excluding the intra-turn pauses to avoid including noisy and redundant signals. For every speaker turn, we extract self-supervised TERA embeddings [16] to obtain a 768 dimensional feature vector. We choose TERA embeddings as it employs a combination of auxiliary tasks to learn the embedding instead of relying on a single task, so it is expected to learn enhanced features from raw speech signals.



**Fig. 3:** Absolute values of CED across age and gender from ADOSMod3

#### 3.2.2. Parameters and implementation details

We use 352 and 64 attention units for the conformer and transformer layers, respectively, while 4 attention heads are employed for both. The full architecture obtained by using a conformer layer followed by a transformer layer results into 2.1M parameters. The model is trained with a binary cross entropy with logits loss function and Adam optimizer with the initial learning rate of  $1e^{-5}$ . There is a provision of early stopping after 10 epochs if no improvement is seen in validation loss, a dropout rate of 0.2 is used for every dropout layer used in the model.

### 3.3. Experimental validation of CED

We carry out an ad-hoc *real/fake* classification experiment to validate CED as a metric for measuring entrainment. For every *real* sample session we synthesize a *fake* sample session by shuffling the speaker turn while maintaining the dyadic conversation sequence. The hypothesis is more entrainment is expected to be observed in *real* sessions as compared *fake* sessions resulting in the *real* sessions having lesser CED. The classification accuracies are reported in Table 2. The classification experiment steps are as follows:

- We calculate CED measure for every consecutive turn pair for the *real* and *fake* sample session.
- We compare the average CED distance from all the turn pairs for the *real* and *fake* session, the sample sessions are correctly classified if CED of *real* session is lesser than *fake* session.
- The experiment is repeated 30 times to eliminate any bias introduced while randomly shuffling the speaker turns.

As baselines, we use three distance measures computed between the extracted turn-level pretrained embeddings: smooth L1 distance [6] (Baseline 1) and two measures introduced in [9], namely, DTWD (Baseline 2), and SCDC (Baseline 3).

**Table 3:** Correlation experiment between CED and clinical scores relevant to ASD (bold figures imply statistical significance,  $p < 0.05$ )  
(*CP: child to psychologist, PC: psychologist to child*)

Clinical scores	Pearson’s correlation			
	CED-PC		CED-CP	
	$\rho$	$p$ -value	$\rho$	$p$ -value
VINELAND ABC	-0.061	0.237	0.012	0.827
VINELAND Social	-0.021	0.345	0.071	0.073
VINELAND Communication	<b>-0.158</b>	<b>0.003</b>	0.043	0.428
CSS	<b>0.222</b>	<b>0.004</b>	0.023	0.672
CSS-SA	<b>0.231</b>	<b>0.012</b>	0.03	0.472
CSS-RRB	0.158	0.055	0.091	0.262

### 3.4. Experimental evaluation

In this experiment, we calculate the correlation between the proposed CED measure and the clinical scores relevant to ASD in Table 3. Since CED is directional in nature, we compute the correlation metric in both the directions *child to psychologist* and *psychologist to child*. We report the Pearson’s correlation coefficient ( $\rho$ ) and also the corresponding  $p$ -value, to test the null hypothesis that there exists no linear association between the proposed measure and the clinical scores. Amongst the clinical scores, *VINELAND* scores are designed to measure adaptive behaviour of individuals, while *VINELAND ABC* stands for Adaptive Behaviour Composite score, *VINELAND social* and *VINELAND communication* are adaptive behavior scores for specific skills of socialization and communication. *CSS* stands for Calibrated Severity Score which reflects the severity of ASD symptoms in individuals. *CSS-SA* and *CSS-RRB* reflects ASD symptoms severity along 2 domains of *Social Affect* and *Restrictive and Repetitive Behaviours*. The details of the clinical scores related to ASD are described in [17, 18, 19].

We also report the absolute values of the proposed CED measure (both directions) for different gender and different age-groups. We partition the dataset across 3 age groups of *Group 1:  $\leq 5yrs$ , Group 2:  $> 5yrs$  &  $\leq 10yrs$ , Group3:  $> 10yrs$*  and for each of the age groups we report the directional CED measure for male and female subgroups in Fig. 3.

## 4. RESULTS AND DISCUSSION

The results reported in Table 2 reveal that achieves better performance in identifying *real* and *fake* sessions with respect to the baseline methods in both Fisher and ADOSmod3 corpus in terms of classification accuracy, which validates the use of CED as a proxy metric for measuring entrainment.

Results in Table 3 show that *VINELAND* communication score is negatively correlated with *psychologist→child* CED with significant statistic, which stands consistent with the definition of CED, since higher CED signifies lower entrainment. *CSS* and *CSS-SA* scores are reported to be positively correlated with CED. It is interesting to note that while *psychologist→child* CED is capturing signals with

meaningful interpretations, no such evidence is reported from *child→psychologist* CED measures. A possible explanation can be since the model is trained with dyadic conversations from adults in Fisher corpus, the model is unable to capture the nuances of interactions involving children which is reflected in these results. It is also worth mentioning while there exists a significant correlation between *CSS*, *CSS-SA* and *psychologist→child* CED, *CSS-RRB* also shows weak evidence of positive correlation with *psychologist→child* CED.

In Table 3, the distributions for absolute values of CED are reported across gender and age-groups. Both directional CED are always seen to have lesser mean values in females as compared to males, which reiterates the claim reported in [20] that women are better at disguising autism symptoms than men. Across age-groups, the experimental results do not show any discernable observation from CED in both directions in male children, however female *psychologist→child* CED is shown to decrease with an increase in age, which also supports the claim presented in [20].

We also investigate the weights of the activations from the cross-encoder attention layer to understand which parts of the speaker turns are emphasized by the attention heads to extract meaningful signals. Attention activation heatmaps from cross-encoder 1 and 2 reported in Fig. 2 show attention layers attend to initial few timeframes from the second speaker turn which supports the claim mentioned in [6] and domain theory that initial and final interpausal units from second and first speaker respectively are a rich source of signals related to entrainment.

## 5. CONCLUSION

In this work we introduce a novel context-aware approach (CED) to measure vocal entrainment in dyadic conversations. We use a combination of convolutional neural networks and transformers to capture both short-term and long-term context, and also employ a cross-subject attention module to learn interpersonal entrainment related information from the other subject in a dyadic conversation. We validate the use of CED as a proxy metric for measuring entrainment by conducting a classification experiment to distinguish between *real* (consistent) and *fake* (inconsistent) interaction sessions. We also study the association between CED and clinically relevant scores related to ASD symptoms by computing the correlation metric. We also report the mean absolute value of directional CED across gender and different age-groups to understand if the entrainment pattern of the children varies across gender or age-group or not. In this work, we use a self-supervised embedding for feature extraction, it will be interesting to see if other context-based pre-trained embeddings yield similar performance in capturing entrainment. We also face difficulties in deploying entrainment embeddings learnt on Fisher for ADOSMod3 dataset and thus we plan to investigate domain-specific entrainment embeddings for understanding behavioral traits.

## 6. REFERENCES

- [1] Alessandro Vinciarelli, Maja Pantic, Hervé Bourlard, and Alex Pentland, “Social signal processing: state-of-the-art and future perspectives of an emerging domain,” in *Proceedings of the 16th ACM international conference on Multimedia*, 2008, pp. 1061–1070.
- [2] Shrikanth Narayanan and Panayiotis G Georgiou, “Behavioral signal processing: Deriving human behavioral informatics from speech and language,” *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1203–1233, 2013.
- [3] Susan E Brennan, “Lexical entrainment in spontaneous dialog,” *Proceedings of ISSD*, vol. 96, pp. 41–44, 1996.
- [4] Rivka Levitan and Julia Bell Hirschberg, “Measuring acoustic-prosodic entrainment with respect to multiple levels and dimensions,” 2011.
- [5] Howard Giles and Peter Powesland, “Accommodation theory,” in *Sociolinguistics*, pp. 232–239. Springer, 1997.
- [6] Md Nasir, Brian Baucom, Craig Bryan, Shrikanth Narayanan, and Panayiotis Georgiou, “Modeling vocal entrainment in conversational speech using deep unsupervised learning,” *IEEE Transactions on Affective Computing*, 2020.
- [7] Chi-Chun Lee, Matthew P. Black, Athanasios Katsamanis, Adam Lammert, Brian Baucom, Andrew Christensen, Panayiotis Georgiou, and Shrikanth S. Narayanan, “Quantification of prosodic entrainment in affective spontaneous spoken interactions of married couples,” in *Proceedings of InterSpeech*, September 2010, pp. 793–796.
- [8] Bo Xiao, Zac E. Imel, David Atkins, Panayiotis Georgiou, and Shrikanth S. Narayanan, “Analyzing speech rate entrainment and its relation to therapist empathy in drug addiction counseling,” in *Proceedings of InterSpeech*, sep 2015.
- [9] Rimita Lahiri, Md Nasir, Manoj Kumar, So Hyun Kim, Somer Bishop, Catherine Lord, and Shrikanth Narayanan, “Interpersonal synchrony across vocal and lexical modalities in interactions involving children with autism spectrum disorder,” *JASA Express Letters*, vol. 2, no. 9, pp. 095202, 2022.
- [10] Michael Natale, “Convergence of mean vocal intensity in dyadic communication as a function of social desirability,” *Journal of Personality and Social Psychology*, vol. 32, no. 5, pp. 790, 1975.
- [11] Nichola Lubold and Heather Pon-Barry, “Acoustic-prosodic entrainment and rapport in collaborative learning dialogues,” in *Proceedings of the 2014 ACM workshop on Multimodal Learning Analytics Workshop and Grand Challenge*, 2014, pp. 5–12.
- [12] Emilie Delaherche, Mohamed Chetouani, Ammar Mahdhaoui, Catherine Saint-Georges, Sylvie Viaux, and David Cohen, “Interpersonal synchrony: A survey of evaluation methods across disciplines,” *IEEE Transactions on Affective Computing*, vol. 3, no. 3, pp. 349–365, 2012.
- [13] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [14] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.
- [15] Christopher Cieri, David Miller, and Kevin Walker, “The fisher corpus: A resource for the next generations of speech-to-text,” in *LREC*, 2004, vol. 4, pp. 69–71.
- [16] Andy T Liu, Shang-Wen Li, and Hung-yi Lee, “Tera: Self-supervised learning of transformer encoder representation for speech,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2351–2366, 2021.
- [17] Somer L Bishop, Marisela Huerta, Katherine Gotham, Karoline Alexandra Havdahl, Andrew Pickles, Amie Duncan, Vanessa Hus Bal, Lisa Croen, and Catherine Lord, “The autism symptom interview, school-age: A brief telephone interview to identify autism spectrum disorders in 5-to-12-year-old children,” *Autism Research*, vol. 10, no. 1, pp. 78–88, 2017.
- [18] Katherine Gotham, Andrew Pickles, and Catherine Lord, “Standardizing ados scores for a measure of severity in autism spectrum disorders,” *Journal of autism and developmental disorders*, vol. 39, no. 5, pp. 693–705, 2009.
- [19] Vanessa Hus, Katherine Gotham, and Catherine Lord, “Standardizing ados domain scores: Separating severity of social affect and restricted and repetitive behaviors,” *Journal of autism and developmental disorders*, vol. 44, no. 10, pp. 2400–2412, 2014.
- [20] Eric Fombonne, “Camouflage and autism,” 2020.