

Automatic Analysis of Typical and Atypical Encoding of Spontaneous Emotion in the Voice of Children

Fabien Ringeval^{1,2}, Erik Marchi^{1,2}, Charline Grossard³,
Jean Xavier³, Mohamed Chetouani³, David Cohen³, Björn Schuller^{1,2,4}

¹Chair of Complex & Intelligent Systems, University of Passau, Germany

²audEERING GmbH, Gilching, Germany

³Institute of Intelligent Systems and Robotics, Université Pierre et Marie Curie, Paris, France

⁴Department of Computing, Imperial College London, UK

fabien.ringeval@uni-passau.de

Abstract

Children with Autism Spectrum Conditions (ASC) present significant difficulties to understand and express emotions. Systems have thus been proposed to provide objective measurements of acoustic features used by children suffering from ASC to encode emotion in speech. However, only a few studies have exploited such systems to compare different groups of children in their ability to express emotions, and even less have focused on the analysis of spontaneous emotion. This contribution aims to fill this white spot in the literature and provides insights by extensive evaluations carried out on a new database of spontaneous speech inducing three emotion categories of valence (positive, neutral, and negative). We evaluate the potential of using an automatic recognition system to differentiate groups of children, i.e., pervasive developmental disorders, pervasive developmental disorders not-otherwise specified, specific language impairments, and typically developing, in their abilities to express spontaneous emotion in a common unconstrained task. Various combinations of feature subsets, i.e., spectral-, source-, and duration-related features, are investigated using Support Vector Machines in a speaker independent approach. Results show that all groups of children can be differentiated directly (diagnosis recognition) and indirectly (emotion recognition) by the proposed system.

Index Terms: affective computing, spontaneous emotions, autism spectrum conditions, language impairments

1. Introduction

The ability to communicate with speech requires the acquisition of codes that link acoustic realisation (e.g., segmental and supra-segmental), to both linguistic [1] and socio-affective related meanings [2, 3]. The acquisition and correct use of such codes, which are supposed to be functional in the early stages of a child's life [4], play an essential role in the inter-subjective development and social interaction abilities of children. As a consequence, most children presenting speech or developmental disorders have limited social interactions, which contributes to social isolation [5].

International classifications differentiate Specific Language Impairment (SLI) from those that are symptomatic of a developmental disorder, e.g., Pervasive Developmental Disorders (PDD) [6]. The former can affect both expressive and receptive language and is defined as a 'pure' language impairment. The latter, PDD, is characterised by severe deficits and perva-

sive impairment in several areas of development such as reciprocal social interactions, communication skills and stereotyped behaviours, interests and activities [7]. Three main developmental disorders have been described [6]: (i) Autism Disorders (AD), with symptoms in all areas that characterise PDD; (ii) Asperger's Syndrome (AS), which does not evince language delay; and (iii) Pervasive Developmental Disorder-Not Otherwise Specified (PDD-NOS), which is characterised by social, communicative and/or stereotypic impairments that are less severe than in AD and appear later in life.

Both clinicians and researchers are facing a huge increase in the prevalence of ASC, resulting from the expansion of the diagnostic criteria, but also from a better awareness of the condition and the acceptance that ASC is a lifelong condition [8]. Recently, systems have been proposed to provide objective measurements of acoustic features used by children suffering from ASC to encode non-verbal information in speech by the use of prosody [9, 10, 11, 12]. Analyses can be then performed indirectly, by assessing the performance of a child on a given task, e.g., producing specific prosodic contours to convey sentence modality [9], or emotion [10, 12]. In this case, the system is tuned for each group of children, e.g., typically developing (TD) and ASC, and performance can be compared between the groups to provide cues regarding the observed atypicalities of ASC. Analyses can also be performed directly, as an automatic diagnosis, by comparing the children's groups in the task [13, 14]. In this case, the system is tuned to search for differences in speech production between each group of children, which can also be a mean to identify the particularities of ASC.

Such systems can be used to help clinicians to improve the diagnosis, but also to develop tools based on information and communication technology, which enable users to access professional support on-line [15]. However, the automatic processing of children's speech is challenging, as they present significant differences compared to the voice of adults [16], and even more when they are affected by ASC. Yet, in spite of these challenges, we do, however, consider the inclusion of emotion recognition technology from speech in ICT-based virtual platform as an important component to enable the development of social imagination and emotional awareness for children suffering from ASC [13, 17, 18, 19]. Recent studies have indeed shown that it is possible to improve emotional skills of ASC children in both emotion perception and production, by providing them interactive tools integrating affective computing technology [15].



Figure 1: Extracts from the book "Frog where are you?" [20] that were used for recording spontaneous emotional speech production from children. We hypothesised that images with *Neutral* valence (left picture, the beginning of the story: the boy is sleeping in his room while his frog is escaping), *Negative* valence (middle picture, the middle of the story: the boy has been captured by a deer while searching for his frog) and *Positive* valence (right picture, the end of the story: the boy finally finds his frog) correlated with emotion production.

Feature set	Spectral	Source	Duration	Total
eGeMAPS	48	48	6	102
IS09	336	48	0	384
IS10	1216	212	154	1582
IS11	2808	272	1288	4368
ComParE	4366	397	1610	6373

Table 1: Distribution of acoustic features in spectral/energy-related, source/excitation-related and duration-related feature sets for different configurations of openSMILE.

1.1. Contribution of this work

The present study focuses on the recognition of spontaneous emotional expressions in the voice of AD, PDD-NOS, SLI and typically developing (TD) children. We investigate the classification performances with expert-based reduced feature sets against large sets of features that include a vast number of spectral-, source- and duration-related features. This study further focuses on the automatic discrimination of typicality between TD children and children suffering from AD, PDD-NOS and SLI. For this purpose, a new database – Child Pathological and Emotional Speech database (CPESD) – is introduced, which will be made available to academic researchers. To the best knowledge of the authors, this is the very first study that performs automatic analysis of spontaneous emotion in children’s speech with AD, PDD-NOS, and SLI.

2. Child Pathological & Emotional Speech Database

We received approval by the Ethical Committee of the Pitié-Salpêtrière Hospital to conduct recruitment and speech recording of children. Consents were obtained from parents or legal caregivers of all participants. Thirty-five monolingual participants with communicative verbal skills were recruited in two university departments of child and adolescent psychiatry located in Paris, France. They consulted for ASC and/or SLI which were diagnosed as AD, PDD-NOS, or SLI, according to DSM IV criteria [6]. Patients were matched for age, sex, academic grades and lexical abilities. Socio-demographic and

Feature set	Negative	Neutral	Positive	All
<i>a. Typicality</i>				
eGeMAPS	83.08	82.60	78.97	83.22
IS09	83.20	80.54	79.69	84.03
IS10	87.09	85.22	85.23	87.59
IS11	89.12	87.46	87.39	89.35
ComParE	88.02	85.12	86.49	86.27
<i>b. Diagnosis</i>				
eGeMAPS	46.15	48.22	44.91	48.23
IS09	50.80	45.78	46.66	51.15
IS10	51.78	47.14	50.09	53.74
IS11	54.74	52.23	51.27	56.38
ComParE	53.24	51.24	50.91	56.17

Table 2: UAR - typicality (2 classes) and diagnosis (4 classes) for each feature set and emotion category.

clinical characteristics of the participants are available in [9]. We also recruited a group of 70 TD children matched for age and sex (1 patient for 2 TD) in elementary schools. A teacher questionnaire was used to exclude children with learning disorders, an history of speech, language, hearing, or general learning problems.

Our main goal was to compare children’s abilities to use prosody to encode pragmatic and affect in speech. A first task was based on the reproduction of intonation contour and was analysed in a previous study [9]. The second task was based on a story telling of a pictured book "Frog where are you?" [20], wherein a little boy tries to find his escaped frog during the night. The task was originally developed to assess language production in a standardised but unconstrained manner. Here, we assume that the child is supposed to produce prosodic cues during the story telling that are correlated to the levels of the emotional valence, which was categorised in three categories by a psychologist: Negative/Neutral/Positive. In total, the pictured book included 15 emotionally negative, 6 emotionally neutral and 5 emotionally positive pictures, cf. Figure 1. Three pictures considered ambivalent because of ambiguous interpretation were excluded.

Group	Valence										
	Negative			Neutral			Positive			All	
	#	%	Duration	#	%	Duration	#	%	Duration	#	Duration
AD	335	35.4 ^{*N}	2.13 ^{*N,T}	137	34.8	2.40 ^{*T}	94	29.8 ^{*N}	2.43 ^{*N,S,T}	566	2.25 ^{*T}
NOS	283	30.1 ^{*A,T}	2.31 ^{*A,S,T}	126	32.2 ^{*T}	2.51 ^{*T}	118	37.7 ^{*A,S,T}	2.08 ^{*A,T}	527	2.31 ^{*S,T}
SLI	530	31.8	2.17 ^{*N,T}	243	35.0	2.30 ^{*T}	184	33.2 ^{*N,T}	2.16 ^{*A,T}	957	2.20 ^{*N,T}
TD	2146	33.8 ^{*N}	2.83 ^{*A,N,S}	970	36.7 ^{*N}	2.89 ^{*A,N,S}	623	29.5 ^{*N,S}	2.78 ^{*A,N,S}	3739	2.84 ^{*A,N,S}

Table 3: Number, relative proportion, and mean duration of utterances per emotion class and group. * = $p < 0.05$: alternative hypothesis is true when comparing data between children groups, i. e., A, N, S and T; AD (A): autism disorders; NOS (N): pervasive developmental disorders not-otherwise specified; SLI (S): specific language impairment; TD (T): typically developing.

Task	eGeMAPS			IS09		IS10			IS11			ComParE		
	Spec.	Sour.	Dur.	Spec.	Sour.	Spec.	Sour.	Dur.	Spec.	Sour.	Dur.	Spec.	Sour.	Dur.
Typ.	83.37	71.11	66.65	81.95	74.77	86.80	76.34	74.91	88.33	80.61	85.72	88.15	72.80	84.41
Diag.	46.26	39.41	34.05	49.97	40.13	53.34	41.99	41.25	56.70	43.51	49.49	56.52	38.91	47.87

Table 4: UAR - typicality (2 classes) and diagnosis (4 classes) for each feature subset and all emotion categories.

Feature set	AD	NOS	SLI	TD
eGeMAPS	36.22	35.62	42.71	44.10
IS09	35.53	39.29	39.34	42.93
IS10	37.69	40.47	41.97	43.10
IS11	38.04	37.94	40.89	44.54
ComParE	35.75	37.04	38.93	42.34

Table 5: UAR - spontaneous emotion recognition (3 classes) for each feature set and groups of children; AD: autism disorders; NOS: pervasive developmental disorders not-otherwise specified; SLI: specific language impairment; TD: typically developing.

We collected nearly 10 hours of recording: 7h38min for TD children, 1h35min for children with AD, 1h12min for children with PDD-NOS, and 1h56min for children with SLI. Recordings were then segmented automatically into groups of breaths, using the energy contour. As many sources of perturbation appeared during the recordings (e. g., false-starts, repetitions or noise from the environment), the obtained speech segments were further manually processed; only utterances that had a complete prosodic contour, i. e., whatever the pronounced words, were kept. Statistics (number, relative proportion, and mean duration) on those utterances are provided for each emotion category, and all, in Table 3. Those data already provide some interesting insights: all TD children produced utterances which are significantly longer than AD, PDD-NOS, and SLI children for all emotion categories ($p < 0.5$, two-tailed t test); we observed the opposite on the constrained task of intonation contour imitation [9]. Moreover, spontaneous speech production of PDD-NOS children focused significantly more on positive emotions compared to all other groups ($p < 0.5$).

3. Experiments

Two main tasks were performed: automatic recognition of typicality (direct analysis), and emotion (indirect analysis). The typicality task concerns the classification of TD children vs. all

others children. Additionally, we performed the classification of each group of children (diagnosis). The emotion task covers the recognition of the three classes of emotional valence, i. e., positive, neutral, and negative. This task was performed either on each group separately, or with models trained on TD children.

3.1. Acoustic features

Acoustic features were automatically extracted from the speech waveform on the utterances by using our open-source openS-MILE feature extractor in its recent 2.1 release [21]. Five different feature sets were investigated: large brute-forced feature sets (IS09, IS10, IS11, and ComParE), which have all been used for paralinguistic information retrieval, and a smaller, expert knowledge based feature set (eGeMAPS). Those feature sets cover spectral-, source- and duration-related feature space with different levels of detail, cf. Table 1. The first four sets, i. e., IS09, IS10, IS11, and ComParE, shows a clear tendency in enlarging the feature space over the years, by including further low-level acoustic descriptors and associated functionals. Recently, this "brute-forcing" approach has been totally revisited, with investigations on a small, expert knowledge based feature set, eGeMAPS [22]. A detailed description and implementation of these feature sets, which is impossible to provide here, is given in [23].

3.2. Setup and evaluation

We used Support Vector Machines (SVMs) for the classification tasks with LIBSVM [24], as they are a well known standard method that can handle both high and low dimensional data. The SVM training has been made with three different kernels: linear, polynomial (3rd order), and Gaussian (γ parameter was set to default value); the complexity parameter was set to default value ($C = 1$). Results are always presented with the best kernel. To ensure speaker independent evaluations, we performed a Leave-One-Speaker-Out (LOSO) cross-validation in all experiments. Because all data sets are unbalanced, we applied upsampling of the under-represented classes in all the evalua-

Group	eGeMAPS			IS09		IS10			IS11			ComParE		
	Spec.	Sour.	Dur.	Spec.	Sour.	Spec.	Sour.	Dur.	Spec.	Sour.	Dur.	Spec.	Sour.	Dur.
AD	34.81	37.44	37.67	37.00	37.41	35.86	34.78	33.60	39.82	37.26	36.84	38.72	33.49	35.75
NOS	38.49	39.09	30.17	38.37	37.07	41.36	36.47	39.15	34.09	36.85	36.98	33.41	34.98	37.04
SLI	43.66	38.46	38.01	38.52	37.12	42.13	38.24	34.90	40.40	37.38	38.67	40.50	37.50	38.93
TD	45.97	41.30	34.63	43.12	35.96	43.66	38.33	38.93	44.35	38.33	42.04	44.23	39.89	42.34

Table 6: UAR - spontaneous emotion recognition (3 classes) for each feature subset and groups of children; AD: autism disorders; NOS: pervasive developmental disorders not-otherwise specified; SLI: specific language impairment; TD: typically developing.

Feature set	AD	NOS	SLI	TD
Spectral + Source	39.97	39.70	43.31	44.73
Source + Duration	37.31	37.27	39.54	42.78
Spectral + Duration	40.64	41.54	39.01	43.21
Spec. + Sour. + Dur.	39.83	42.06	39.06	43.22
Best group	39.82	41.36	43.66	45.97

Table 7: UAR - spontaneous emotion recognition with 3 classes, combination of best feature subsets; AD: autism disorders; NOS: pervasive developmental disorders not-otherwise specified; SLI: specific language impairment; TD: typically developing.

Feature set	AD	NOS	SLI	TD
Spectral	35.81	41.65	35.76	45.97
Source	39.33	38.43	33.94	41.30
Duration	34.16	37.78	38.01	42.34
Spectral + Source	38.02	43.47	33.16	44.73
Source + Duration	36.34	39.16	37.47	42.78
Spectral + Duration	34.72	40.54	38.73	43.21
Spec. + Sour. + Dur.	35.10	40.72	38.20	43.22
IS11	38.30	40.82	39.59	44.54

Table 8: UAR - spontaneous emotion recognition with 3 classes and model training on CTR; AD: autism disorders; NOS: pervasive developmental disorders not-otherwise specified; SLI: specific language impairment; TD: typically developing.

tion experiments. For the same reason, we used the unweighted average recall (UAR) of the classes as scoring metric. Standardisation of the features, i. e., feature values are normalised to zero-mean and unit standard deviation, was performed for each speakers for all emotion recognition tasks. Whereas for both typicality and diagnosis tasks, we standardised the features of all speakers with the on-line approach, i. e., mean and standard-deviation were computed on the training partition and applied on both training and testing data.

4. Results

4.1. Typicality & Diagnosis

Results obtained on typicality and diagnosis are given for each emotion class and each feature set in Table 2. One may note that all obtained performance are far above the chance level, in agreement with [12], despite being lower than on the constrained task, i. e., intonation contour imitation [25]. In order to gain further insights, we performed automatic recognition of

typicality and diagnosis with each different feature subset, i. e., spectral-, source-, and duration-related features. Results show that, spectral-related features are the most contributing in the two tasks, and can perform even better when taken alone for diagnosis, cf. Table 4.

4.2. Emotion

Results obtained on the automatic recognition of the emotional valence are given in Table 5. Our hypothesis that the spontaneous description of the pictured book will be correlated with the emotional valence depicted in the images is validated by the experiments, because the system performs significantly better than chance for all groups of children. Obviously, TD children obtained the best performance, and all others children obtained a significantly lower performance ($p < 0.5$). A detailed analysis of each feature subset shows that, (i) spectral-related features provide once more the best performance for all groups of children, and (ii) the minimalistic feature set, i. e., eGeMAPS, performed remarkably well on all groups, especially for source-related features, cf. Table 6. Results obtained with different combinations of the best feature subsets show that, the performance was improved further for the PDDs, by combining spectral- and duration-related features for AD, and all feature subsets for PDD-NOS, cf. Table 7.

Finally, in order to investigate how the models obtained on TD children could generalise on the others groups of children, we performed a mismatched evaluation, by training models on TD and testing on AD, PDD-NOS, and SLI. Results show that, there are some specific associations between feature space and pathology; models obtained from TD children generalised best on AD with source-related features, PDD-NOS with spectral-related features, and SLI with duration-related features. Moreover, PDD-NOS children obtained systematically the best performance for all combinations of the feature subsets, cf. Table 8.

5. Conclusions

A new speech database of spontaneous emotions produced by AD, PDD-NOS, SLI and TD French speaking children has been introduced: CPESD. Extensive experiments have been performed on this database, showing that all groups of children can be differentiated directly (typicality, diagnosis) and indirectly (emotion) with an automatic recognition system.

6. Acknowledgements

The research leading to these results has been funded by the European Community’s Seventh Framework Programme through the ERC Starting Grant No. 338164 (iHEARu), and the European’s Union’s Horizon 2020 Programme through the Research Innovative Action #688835 (DE-ENIGMA).

7. References

- [1] D. Crystal, "The linguistic status of prosodic and paralinguistic features," *Proc. of the University of Newcastle-upon Tyne Philosophical Society*, vol. 1, pp. 93–108, 1966.
- [2] ———, "Prosodic and paralinguistic correlates of social categories," in *Social anthropology*, E. Ardener, Ed. London: Tavistock, 1971, pp. 185–206.
- [3] E. Marchi, Y. Zhang, F. Eyben, F. Ringeval, and B. Schuller, "Autism and Speech, Language, and Emotion – a Survey," in *Evaluating the Role of Speech Technology in Medical Case Management*, H. Patil and M. Kulshreshtha, Eds. Berlin: De Gruyter, 2015, invited contribution, to appear.
- [4] P. K. Kuhl, "Early language acquisition: cracking the speech code," *Nature Reviews Neuroscience*, vol. 5, pp. 831–843, November 2004.
- [5] E. Müller, A. Schuler, and G. B. Yates, "Social challenges and supports from the perspective of individuals with Asperger syndrome and other autism spectrum disabilities," *Autism*, vol. 12, no. 2, pp. 173–190, March 2008.
- [6] "Diagnostic and Statistical Manual of mental disorders (4th Ed.)." Washington, DC: American Psychiatric Association, 1994.
- [7] I. Rapin and D. A. Allen, "Developmental language: nosological consideration," in *Neuropsychology of Language, Reading and Spelling*, V. Kvik, Ed. Washington, DC: New York: Academic Press, 1983.
- [8] J. L. Matson and A. M. Kozlowski, "The increasing prevalence of autism spectrum disorders," *Research in Autism Spectrum Conditions*, vol. 5, no. 1, pp. 418–425, January–March 2011.
- [9] F. Ringeval, J. Demouy, G. Szaszák, M. Chetouani, L. Robel, J. Xavier, D. Cohen, and M. Plaza, "Automatic intonation recognition for prosodic assessment of language impaired children," *IEEE Transactions on Audio, Speech & Language Processing*, vol. 19, no. 5, pp. 1328–1342, July 2011.
- [10] E. Marchi, B. Schuller, A. Batliner, S. Fridenzon, S. Tal, and O. Golan, "Emotion in the Speech of Children with Autism Spectrum Conditions: Prosody and Everything Else," in *Proc. of WOCCI 2012*, ISCA. Portland, OR: ISCA, September 2012.
- [11] D. Bone, C.-C. Lee, M. P. Black, M. E. Williams, S. Lee, P. Levitt, and S. Narayanan, "The psychologist as an interlocutor in autism spectrum disorder assessment: Insights from a study of spontaneous prosody," *Journal of Speech, Language, and Hearing Research*, vol. 57, no. 4, pp. 1162–1177, August 2014.
- [12] E. Marchi, B. Schuller, S. Baron-Cohen, O. Golan, S. Bölte, P. Arora, and R. Hüb-Umbach, "Typicality and Emotion in the Voice of Children with Autism Spectrum Condition: Evidence Across Three Languages," in *Proc. of INTERSPEECH 2015*, ISCA. Dresden, Germany: ISCA, September 2015, pp. 115–119.
- [13] E. Marchi, F. Ringeval, and B. Schuller, "Voice-enabled assistive robots for handling autism spectrum conditions: an examination of the role of prosody," in *Speech and Automata in Health Care (Speech Technology and Text Mining in Medicine and Healthcare)*, A. Neustein, Ed. Boston/Berlin/Munich: De Gruyter, 2014, pp. 207–236.
- [14] E. Marchi, A. Batliner, B. Schuller, S. Fridenzon, S. Tal, and O. Golan, "Speech, Emotion, Age, Language, Task, and Typicality: Trying to Disentangle Performance and Feature Relevance," in *Proc. WS³P 2012, SocialCom 2012*, ASE/IEEE. Amsterdam, The Netherlands: IEEE, September 2012, 8 pages.
- [15] B. Schuller *et al.*, "Recent developments and results of ASC-Inclusion: An Integrated Internet-Based Environment for Social Inclusion of Children with Autism Spectrum Conditions," in *Proc. of IDGEI 2015*. Atlanta, GA: ACM, March 2015.
- [16] A. Potamianos and S. Narayanan, "A review of the acoustic and linguistic properties of childrens speech," in *Proc. of the IEEE 9th Workshop on Multimedia Signal Processing*, October 2007, pp. 22–25.
- [17] E. Mower, M. P. Black, E. Flores, M. Williams, and S. Narayanan, "Rachel: Design of an emotionally targeted interactive agent for children with autism," in *Proc. of ICMCS/ICME*, 2011, pp. 1–6.
- [18] O. Golan, S. Baron-Cohen, J. J. Hill, and M. D. Rutherford, "The 'reading the mind in the voice' test-revised: A study of complex emotion recognition in adults with and without autism spectrum conditions," *Journal of Autism and Developmental Disorders*, vol. 37, pp. 1096–1106, 2007.
- [19] O. Golan and S. Baron-Cohen, "Systemizing empathy: Teaching adults with asperger syndrome or high-functioning autism to recognize complex emotions using interactive multimedia," *Development and Psychopathology*, vol. 18, no. 02, pp. 591–617, 2006.
- [20] M. Mayer, *Frog where are you?* New York: Dial Books for young readers, 1969.
- [21] F. Eyben, F. Weninger, F. Groß, and B. Schuller, "Recent developments in openSMILE, the Munich open-source multimedia feature extractor," in *Proc. of the 21st ACM International Conference on Multimedia (ACM MM)*, Barcelona, Spain, October 2013, pp. 835–838.
- [22] F. Eyben, K. Scherer, B. Schuller, J. Sundberg, E. André, C. Busso, L. Devillers, J. Epps, P. Laukka, S. Narayanan, and K. Truong, "The Geneva Minimalistic Acoustic Parameter Set (GeMAPS) for Voice Research and Affective Computing," *IEEE Transactions on Affective Computing*, 2015, 14 pages, to appear.
- [23] F. Eyben, "Real-time speech and music classification by large audio feature space extraction," Ph.D. dissertation, Technische Universität München, 2014, springer.
- [24] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, no. 3, April 2010, Article No. 27.
- [25] B. Schuller *et al.*, "The INTERSPEECH 2013 computational paralinguistics challenge: Social signals, conflict, emotion, autism," in *Proc. of INTERSPEECH 2013*, Lyon, France, August 2013, pp. 148–152.