

Emotion, Age, and Gender Classification in Children’s Speech by Humans and Machines

Heysem Kaya^{a,*}, Albert Ali Salah^b, Alexey Karpov^{c,d}, Olga Frolova^e, Aleksey Grigorev^e, Elena Lyakso^e

^a*Department of Computer Engineering, Namık Kemal University, Corlu, Tekirdag, Turkey*

^b*Department of Computer Engineering, Bogazici University, Istanbul, Turkey*

^c*Speech and Multimodal Interfaces Lab., St. Petersburg Institute for Informatics and Automation of the Russian Academy of Sciences, St. Petersburg, Russia*

^d*Department of Speech Information Systems, ITMO University, St. Petersburg, Russia*

^e*Child Speech Research Group, St. Petersburg State University, St. Petersburg, Russia*

Abstract

In this article, we present the first child emotional speech corpus in Russian, called “EmoChildRu”, collected from 3-7 years old children. The base corpus includes over 20K recordings (approx. 30 hours), collected from 120 children. Audio recordings are carried out in three controlled settings by creating different emotional states for children: playing with a standard set of toys; repetition of words from a toy-parrot in a game store setting; watching a cartoon and retelling of the story, respectively. This corpus is designed to study the reflection of the emotional state in the characteristics of voice and speech and for studies of the formation of emotional states in ontogenesis. A portion of the corpus is annotated for three emotional states (comfort, discomfort, neutral). Additional data include the results of the adult listeners’ analysis of child speech, questionnaires, as well as annotation for gender and age in months. We also provide several baselines, comparing human and machine estimation on this corpus for prediction of age, gender and comfort state. While in age estimation, the acoustics-based automatic systems show higher performance, they do not

☆This is the unedited author proof of the accepted article.

*Corresponding author

Email addresses: hkaya@nku.edu.tr (Heysem Kaya), salah@boun.edu.tr (Albert Ali Salah), karpov@iias.spb.su (Alexey Karpov), olchel@yandex.ru (Olga Frolova), a.s.grigoriev89@gmail.com (Aleksey Grigorev), lyakso@gmail.com (Elena Lyakso)

reach human perception levels in comfort state and gender classification. The comparative results indicate the importance and necessity of developing further linguistic models for discrimination.

Keywords: emotional child speech, perception experiments, spectrographic analysis, emotional states, age recognition, gender recognition, computational paralinguistics

1. Introduction and Related Work

Speech based communication contains both linguistic and paralinguistic information. The latter is particularly important in specifying factors of behavioral and functional status, and especially emotional states. For children's communications, self-reporting is not very reliable as a measure, and assessment of emotional speech becomes particularly valuable. There are two main approaches or the study of emotional speech. One approach focuses on the psychophysiological aspects of emotions, which can include studies of brain activity data (Lindquist et al., 2012; Watson et al., 2014), and cross-cultural investigation of emotional states in speech (Lyakso and Frolova, 2015; Rigoulot et al., 2013; Jürgens et al., 2013; Laukka et al., 2013). The second approach is focused on the physical speech signal and its analysis. Hence, it is geared towards software applications for human-computer interaction, such as automatic speech recognition (Fringi et al., 2015; Liao et al., 2015; Guo et al., 2015) and speech synthesis (Govender et al., 2015).

Adults perceive emotional states of infants in their vocalizations from the first months onwards (Lyakso and Frolova, 2015). For instance discomfort and comfort conditions of three months old infants are recognizable by people, but also via spectrographic analysis, which reveals differences in the pitch values and the duration of vocalizations. Crying and squeals of joy are indicative of emotional states, but acoustic features are not always significantly different. With child's increasing age, lexical information acquires more discriminative power in the recognition of emotional states (Yildirim et al., 2011).

24 It is well known that acoustic and linguistic characteristics of child speech
25 are essentially different from those of adult speech. The child speech is char-
26 acterized by a higher pitch value, formant frequencies and specific indistinct
27 articulation with respect to the adult speech. Recognition of child’s speech can
28 be challenging. It was shown that adult Russians recognize between half and
29 three quarters of 4-5 years old children’s words and phrases in calm and spon-
30 taneous conditions (Lyakso et al., 2006). The paralinguistic aspects, however,
31 require more research, both from a human perceptual perspective, and from
32 an automated speech processing perspective. This paper aims to address these
33 points.

34 The first requirement for studying children’s emotional speech is the prepara-
35 tion of an adequate corpus (Ververidis and Kotropoulos, 2006). Creation of
36 such a corpus is more difficult than the collection of emotional speech corpora of
37 adults. In the case of adults, actors are often involved to portray the necessary
38 emotional conditions (Engberg and Hansen, 1996; Burkhardt et al., 2005; Kaya
39 et al., 2014; Lyakso and Frolova, 2015), or records of patients from a psychi-
40 atry clinic are used. Such approaches are not easily used for children. It is
41 necessary to model communicative tasks in which the child is not conscious of
42 being recorded to produce veridical emotional reactions. The creation of the
43 corpus should be based on a verified and clear method of obtaining spontaneous
44 speech manifestations of certain emotional reactions. By nature, collection of
45 child emotional speech data should be under natural conditions that are not-
46 controlled, not-induced (i.e. “spontaneous”).

47 At present there are a few spontaneous or emotional child speech databases
48 available for the child speech research community. These include emotional and
49 spontaneous corpora for Mexican Spanish (7-13 years old) (Pérez-Espinosa et al.,
50 2011), British English (4-14 years old) (Batliner et al., 2005), and German (10-13
51 years old) (Batliner et al., 2005, 2008). The SpontIt corpus is spontaneous child
52 speech in Italian (8-12 years old) (Gerosa et al., 2007), and the NICE corpus is
53 spontaneous child speech in Swedish, possibly emotional, but without emotion
54 annotations (8-15 years old) (Bell et al., 2005). Recently, we have collected

55 the first emotional child speech corpus in Russian, called “EmoChildRu”, and
56 reported initial results (Lyakso et al., 2015). The present work greatly extends
57 the scope of investigation on this corpus, doubling the annotated data, and
58 providing age and gender estimation baselines for both machine classification
59 and human perceptual tests.

60 The rest of the article is structured as follows: Section 2 introduces the
61 Emotional Child Russian Speech Corpus “EmoChildRu”, including the record-
62 ing setup and speech data analysis. Section 3 describes two separate human
63 perception experiments, one on the recognition of emotional states and another
64 for prediction of child’s age and gender by listeners, respectively. Section 4
65 presents baseline automatic classification systems for paralinguistic analysis,
66 and reports extensive experimental results. Section 5 provides a discussion of
67 the findings and conclusions.

68 **2. Emotional Child Russian Speech Corpus**

69 “EmoChildRu” is the first database containing emotional speech material
70 from 3–7 year old Russian children. Three emotional states (discomfort, com-
71 fort, neutral) are used in the database. It is important to note that the “discom-
72 fort” state encapsulates a number of basic emotions, such as “sadness,” “fear,”
73 and “anger,” but these emotional statements are not expressed strongly. It is
74 not ethical to induce natural fear or anger in 3–7 year old children for the pur-
75 poses of such a study. All procedures were approved by the Health and Human
76 Research Ethics Committee (HHS, IRB 00003875, St. Petersburg State Uni-
77 versity) and written informed consent was obtained from parents of the child
78 participant.

79 All children in the database were born (and lived) in the city of St. Peters-
80 burg, with parents who were also born in St. Petersburg, or have been living
81 there for more than 10 years. The whole collection includes 20.340 utterances
82 (more than 30 hours of speech). Recordings were made at home, in laboratory
83 and at kindergarten. The three different recording conditions are playing with

84 a standard set of toys, repetition of words from a toy-parrot in a game store
85 setting, and watching a Russian cartoon called “Masha and bear” from iPad
86 and the retelling of the story, respectively. Each experiment had a duration of
87 2 minutes (containing multiple utterances). Every record is accompanied by a
88 protocol, which describes the recording conditions, and video recording of child’s
89 behavior in parallel. The speech materials are grouped based on the protocol
90 and on the recording situation.

91 Model situations for provoking the child’s emotional states were selected
92 based on our previous experience - the supervision of children in various forms
93 of interaction with adults (experimenters), and taking into account the response
94 of children aged 4-7 years (Lyakso et al., 2010b). Using heart rate data, the
95 child’s emotional state was estimated during the preliminary experiments (which
96 makes it possible to compare neutral and arousal states), and an additional video
97 analysis was performed by experts. They annotated child behaviors, as well as
98 facial expressions during these preliminary experiments.

99 The speech recordings were made with a “Marantz PMD660” digital recorder
100 and with a single “SENNHEISER e835S” external microphone. The speech
101 sounds were analyzed by experts, using Syntrillium’s “Cool Edit Pro” sound
102 editor. Speech files are stored in Windows PCM format, 22.050 Hz, 16 bits
103 per sample. Stressed vowels were selected from stressed words of all phrases.
104 Pitch and the vowel duration, as well as phrase prosody, were automatically
105 calculated, based on the algorithms implemented in “Cool Edit Pro” sound ed-
106 itor. The waveform view was used for calculation of duration, and the spectral
107 view was used to measure the pitch control for prosody. The corpus and the
108 accompanying software package include the database, as well as a shell compo-
109 nent to navigate, browse and search information in the database. So far, about
110 25% of the data are annotated for emotional states. The child’s (ground truth)
111 emotional state was determined based on the recording setting and by analysis
112 of the video clips by five experts having professional experience of working with
113 children and child speech. The database contains additional information about
114 the child’s psychophysiological condition before and after speech recording. Di-

115 alogue descriptions, speech developmental and cognitive scale data are included
116 in the database, whenever available. We will focus on the acoustic information
117 only for the purposes of this paper.

118 A software tool was developed in Microsoft Visual C# for enabling the ex-
119 perts to work with the “EmoChildRu” corpus under Windows OS. This software
120 also allows choosing speech material using a query interface, along dimensions
121 such as the type of emotional state, child age, and gender.

122 We performed a qualitative analysis of children’s words reflecting different
123 emotional states, using words from 4-years olds (4.800 words from 39 children),
124 5-years olds (9.030 words from 55 children), 6-years olds (4.150 words from 26
125 children) and 7-years olds (846 words from 18 children). 4-year old children
126 express themselves by antonyms (*yes – no; it is terrible – well; I’m afraid; I am*
127 *glad; good – bad*). At the age of 7, the word range significantly expands (e.g.
128 *very angry; angry; terrible; bad; not so good* for discomfort and *like; good; like*
129 *more; like most; love; immense; wonderful; splendid* for comfort). In line with
130 a recent study on a portion of this corpus (Lyakso et al., 2016), the number of
131 words 7 year old children use to reflect a discomfort state ($n = 14 \pm 8$ words)
132 is found higher than the number of words that reflect the comfort state ($n =$
133 10 ± 6 words).

134 2.1. Dataset Used for Machine Learning Experiments

135 For automatic recognition experiments, we used a subset of the corpus (1.116
136 child speech utterances), where all speech files have 1 to 5 seconds of speech
137 signal and all five child speech experts agree on the emotion annotation. Note
138 that the annotation is done from the audio-visual data, including linguistic
139 information, though for automatic classification, only acoustic features are used.

140 The subset contains data from 113 children, and the number of speech files
141 per child ranges from 1 to 78 (mean=9.9), and this imbalance makes auto-
142 matic recognition more difficult. There are 54 boys and 59 girls in the dataset
143 (mean±std age is 5.3 ± 1.1 years). We use a different subset for the machine
144 experiments, because the number of instances needed for a machine classifier is

145 much higher than what can feasibly be set in a human perception study. Also,
 146 we split the data into training and test sets, which is not required in the human
 147 perception study. Class distribution for the three classification tasks over the
 148 speaker-independent training and test partitions are given in Tables 1 and 2.

Table 1: Distribution of emotion and gender classes. #Inst. denotes number of utterances.

Set	Gender		Emotion		
	#M/#Inst	#F /#Inst.	Comfort	Discomfort	Neutral
Train	36 / 288	34 / 353	245	111	285
Test	18 / 209	25 / 266	189	94	192
Total	54 / 497	59 / 619	434	205	477

Table 2: Distribution of age group classes in train and test partitions.

Set	#Inst.	3-4 Years	5 Years	6-7 Years
Train	641	130	168	343
Test	475	111	115	249
Total	1116	241	283	592

149 2.2. Dataset Used for Human Perception Experiments

150 The dataset used for human perception experiments is a subset of the one
 151 used for the automatic recognition experiments. This is because we need a large
 152 number of instances to train an automatic recognizer as opposed to a naturally
 153 trained human.

154 30 children, aged from 3 to 7 years were selected for human perception study,
 155 and three test sequences were formed from the speech material of each child.
 156 These sequences were arranged such that they equally represent the three age
 157 groups (3–4 year old, 5 year old, 6–7 year old children, respectively) and that
 158 every test sequence includes 10 phrases uttered by children in a comfortable
 159 emotional state, 10 phrases in a discomfort state and 10 phrases in a neutral
 160 (calm) state. In total, we used 90 sequences for testing.

161 **3. Human Perceptual Experiments**

162 This section reports two human perceptual experiments to provide insight on
163 the nature of the EmoChildRu Database. Listeners were Pediatric University
164 Students 300 adults (age: 18.8 ± 2.2 years, median 18 years; 61 male, 239
165 female; 219 with the experience of interaction with children). Child interaction
166 experience implies the presence of children in the family – younger brothers and
167 sisters, communication with children of friends and relatives. Data about the
168 listeners with experience and without experience of interaction with the children
169 are presented together, as significant differences in the recognition of children
170 were not found between these groups. The presentation of test sequences was
171 carried out in an open field for 10-people groups (the listener location from
172 the source of sounds had no effect on the recognition result). Each signal in
173 the test (phrase) was presented one, the duration of pauses between the signals
174 was 7 seconds, which allowed the listeners to fill in the forms with requested
175 information. The 7-second interval is chosen experimentally.

176 *3.1. Human Perception and Spectrographic Analysis of Emotional Speech*

177 The aim of the first study is to reveal how humans (Russian native speakers)
178 can recognize emotional states in children via speech. Each of the test sequences
179 was presented to 100 adults for perceptual analysis, and the ratio of listeners
180 correctly recognizing the emotional state on the base of speech samples (percep-
181 tion rate) was calculated. Confusion matrices for perception experiments were
182 prepared.

183 In terms of spectrographic analysis, we analyzed and compared pitch values,
184 max and min values of pitch, pitch range, energy and duration. The vowel
185 duration and pitch values were determined based on the analysis of the dynamic
186 spectrogram. Spectral analysis was performed by fast Fourier transformation
187 weighted using a Hamming window, with a window length of 512 samples. To
188 consider word stress development, the vowel duration and its stationary part
189 duration were compared in the stressed versus the unstressed vowels, as well as
190 the pitch and formant values in the stationary parts.

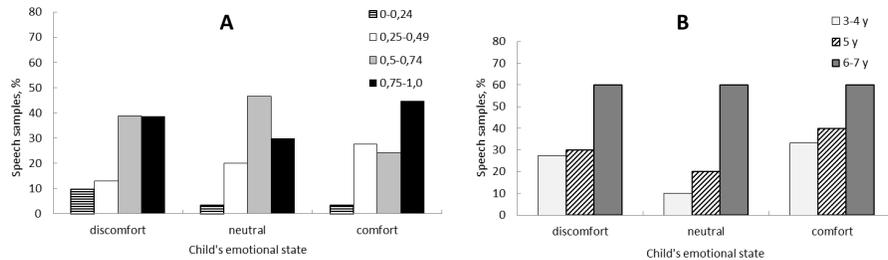


Figure 1: Percentages of emotional child speech samples perceived by the listeners: A) correctly recognized states as ‘discomfort’, ‘neutral’ and ‘comfort’ with the perception rate of 0-0.24 (horizontal hatch), with the rate of 0.25-0.49 (white color bar), with the rate of 0.5-0.74 (light gray) and with the rate of 0.75-1 (black); B) correct recognition with the rate of 0.75-1.0 for different age groups: 3-4 years (light gray), 5 years (sloping hatch), 6-7 years old (grey).

191 Stressed and unstressed vowels in the words are defined initially by experts
 192 via auditory perception of the word meaning in the Russian language. A sec-
 193 ondary (instrumental) spectrographic analysis was performed on stressed vowel
 194 duration and pitch, because a former study found that the stressed vowels in the
 195 child speech are characterized by higher values of the pitch and duration (Lyakso
 196 and Gromova, 2005), unlike the stressed vowels in Russian adult speech that
 197 only exhibit longer duration compared to the unstressed.

198 The average recognition accuracy was 56% for the 3-4 year old group, 66%
 199 for the 5 year old group, and 73% for 6-7 year old group. The comparison of per-
 200 ception rate for emotional speech samples showed that the amount of samples
 201 attributed by most of listeners (0.75-1.0) as the state of comfort was higher than
 202 the amount of signals attributed as neutral and discomfort (Figure 1-A). The
 203 agreement among the listeners in determining neutral signals and discomfort
 204 was lower: majority of neutral signals were recognized with a ratio of 0.5-0.74
 205 (Figure 1-A). The state of comfort is more clearly reflected in the child speech
 206 compared to the state of discomfort. Often 3-6 year old children spoke excit-
 207 edly, smiled and laughed in comfortable emotional state during recording. At
 208 the same time, children of pre-school age do not express intense discomfort,
 209 only a slight discomfort is considered admissible in recording situation in the

210 presence of an unfamiliar adult. The slight discomfort was recorded when child
 211 spoke about something unpleasant, manifested disgust or anger. Consequently,
 212 listeners recognized slight discomfort of child with a smaller ratio compared
 213 to comfort. The neutral state is recognized by listeners with ratio of 0.5-0.74,
 214 because they attributed the part of neutral speech samples to comfort and dis-
 215 comfort samples.

216 Amount of emotional samples correctly recognized by most of listeners (i.e.
 217 perception rate of 0.75-1.0) increased with child age. The emotion perception
 218 rate for the 6-7 year old children was higher than that of 3-4 year old children,
 219 as shown in Figure 1-B ($p < 0.01$, Mann-Whitney U test).

220 Errors in emotional state classification were associated primarily with the
 221 allocation of discomfort and comfort speech samples to the neutral state class.
 222 Speech samples that reflected a neutral state were classified more often as dis-
 223 comfort rather than comfort. The corresponding confusion matrices are pre-
 224 sented in Table 3.

Table 3: Confusion matrices for emotion recognition by humans with age group breakdown.

	Child's age								
	3-4 years			5 years			6-7 years		
State	disc	neut	comf	disc	neut	comf	disc	neut	comf
discomfort	64	22	14	68	23	9	69	25	6
neutral	23	56	21	24	65	11	18	70	12
comfort	13	33	54	9	26	65	4	18	78

225 Spectrographic analysis also revealed that speech samples interpreted by
 226 humans as discomfort, neutral and comfort are characterized by specific sets
 227 of acoustic features (Figure 2). Discomfort speech samples are characterized
 228 by phrase duration lengthening; highest pitch values (comparatively) in phrase
 229 and in stressed vowels selected from stressed words in the phrase ($p < 0.05$ –
 230 vs. neutral state, Mann-Whitney U test); an increase of minimum pitch value
 231 in phrases; an increase of the pitch range in stressed vowels of stressed words

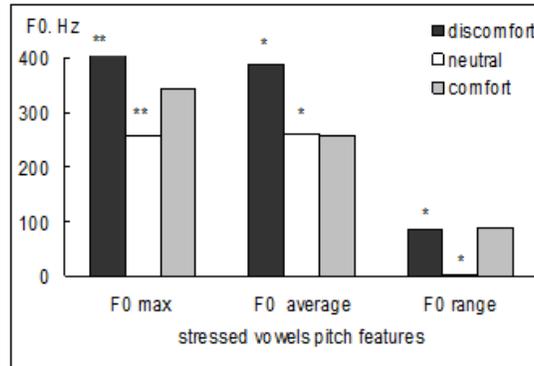


Figure 2: Pitch features of stressed vowels (medians)– F0 max: maximum pitch value, Hz; F0 average: vowel’s average pitch value, Hz; F0 range: vowel’s pitch range value (F0max - F0min); * - $p < 0.05$, ** - $p < 0.01$ Mann-Whitney U test.

232 in the phrase ($p < 0.05$ – vs. neutral state); falling pitch contour of stressed
 233 vowels of stressed words in the phrase. Comfort speech phrases have short
 234 duration, together with long stressed vowel duration; pitch values of phrases
 235 are increased, but less so compared to discomfort samples; pitch range in the
 236 stressed vowels is similar to discomfort samples; pitch contours of stressed vowels
 237 are rising. Neutral speech samples are characterized by lowest values of vowel
 238 duration, stressed vowels’ pitch and pitch range (Figure 2). Flat pitch contours
 239 are observed for 3–4 and 6–7 years old children (Table 4).

240 Most speech samples that were correctly recognized by humans in the ex-
 241 periment have a complex shape of phrase pitch contours ($> 70\%$ samples).
 242 The analysis of features of all stressed vowels from stressed words revealed that
 243 discomfort speech samples have mainly a falling shape, while comfort speech
 244 samples have a rising shape.

245 Almost all neutral speech samples have flat, falling and bell-shaped pitch con-
 246 tours, and the first two patterns are the most common. U-shaped pitch contour
 247 is revealed in comfort speech samples only. The variety of pitch contour shapes
 248 in stressed vowels increases by 6-7 years, compared to younger children. With
 249 increasing age, the duration of phrases increases, and the duration of stressed

Table 4: Distribution of pitch contour shapes for speech samples with correctly recognized emotional states.

		Pitch contour shape, %				
State	Age (y)	flat	rising	falling	U-shaped	bell-shaped
discomfort	3-4	0	33	67	0	0
	5	0	0	100	0	0
	6-7	33	0	67	0	0
neutral	3-4	100	0	0	0	0
	5	0	0	100	0	0
	6-7	67	0	16.5	0	16.5
comfort	3-4	0	67	0	33	0
	5	0	75	25	0	0
	6-7	17	50	0	33	0

250 vowels and pitch values decreases. The differences in the acoustic characteristics
 251 of speech samples correctly recognized as discomfort, neutral and comfort are
 252 more expressed at the age of 3-5 years. Correctly recognized speech samples
 253 of 6-7 years old children do not differ significantly in acoustic features. Adult
 254 listeners mostly rely on the meaning of the phrase. Analysis of speech sam-
 255 ples correctly recognized by listeners revealed that detection of word meaning
 256 improved with increase of child’s age up to 100% at the age of 6-7 years.

257 3.2. Estimation/Perception of Child’s Age/Gender by Humans

258 The aim of the second study was to investigate the possibility of child’s
 259 age and gender recognition by listeners. Three test sequences were formed
 260 from the speech material of 45 children. The sequences contain speech data
 261 uttered by children in a discomfortable emotional state (Test 1 – discomfort),
 262 in a neutral/calm state (Test 2 – neutral) and in a comfortable state (Test 3 –
 263 comfort). Every test sequence includes 30 phrases: 10 speech samples for each
 264 age group (3-4- years, 5 years, and 6-7 years; five speech samples per gender).
 265 Thus, each test sequence includes 15 speech samples uttered by boys and 15

266 samples by girls. For testing, we used 90 sequences in total; each speech signal
267 was included in the test sequence only once; the time interval between the signals
268 was seven seconds.

269 It should be noted that Russian is a gender-dependent language and most of
270 the verbs in the past tense form and some adjectives have both masculine and
271 feminine gender word-forms. We have excluded such phrases from the test set
272 for the human perception experiments, so that Russian speaking listeners could
273 not easily predict child’s gender using linguistic knowledge. In the gender and
274 age prediction tasks, each test sequence was evaluated by 100 adult listeners (a
275 total of 300 listeners), who were asked to select the gender (male or female) and
276 age group (3-4, 5, or 6-7 years) of each child in a questionnaire.

277 For gender prediction, the average recognition accuracies were 66%, 64%
278 and 71% for speech samples uttered by children in the discomfort state (Test-
279 1), neutral state (Test-2), and comfort state (Test-3), respectively. Most of
280 listeners (0.75-1.0) correctly recognized the gender of child on the base of child
281 speech samples reflected neutral, discomfort, and especially comfort state (see
282 Figure 3). The amount of comfort samples on which most of listeners (0.75-1.0)
283 correctly recognized child’s gender was higher than the amount of discomfort
284 and neutral samples. The agreement among the listeners in determining child’s
285 gender on the base of neutral signals was the lowest. Child gender was recognized
286 with ratio 0-0.5 in more than 30% of neutral samples. We can conclude that
287 the emotional child speech includes more acoustic and linguistic information for
288 gender recognition as opposed to neutral speech.

289 Percentages of correct and incorrect gender estimation are reported individ-
290 ually for comfort, neutral and discomfort states with confusion matrices in Ta-
291 ble 5. The least number of errors was made for male speakers at the discomfort
292 state, and the most error-prone classification was female speaker recognition at
293 the comfort state. Human recognition of male speakers was better than female
294 speakers for all emotional states.

295 In the age prediction task, the average recognition accuracy was 50%, 52%,
296 and 51% of speech samples uttered by children in the discomfort state, neu-

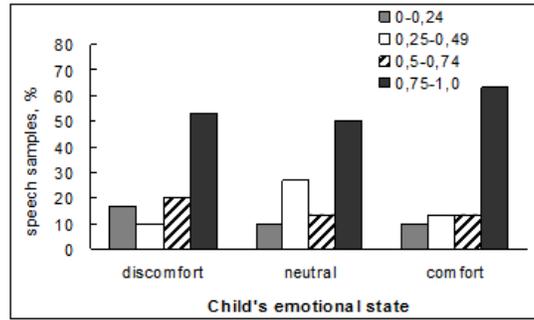


Figure 3: Percentages of emotional child speech samples perceived by the listeners: correctly recognized gender with the perception rate of 0-0.24 (gray), with the rate of 0.25-0.49 (white color), with the rate of 0.5-0.74 (sloping hatch) and with the rate of 0.75-1 (black)

Table 5: Confusion matrices for gender prediction in three emotional states.

	Emotional state					
	discomfort		neutral		com fort	
Gender	male	female	male	female	male	female
male	80	20	68	32	76	24
female	48	52	40	60	34	66

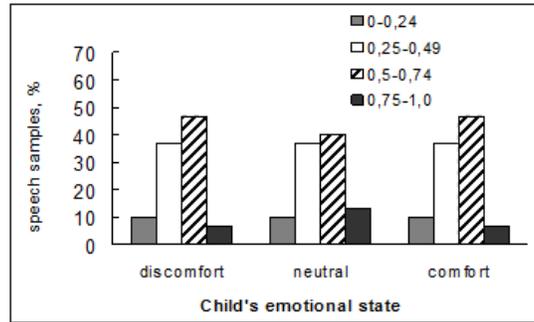


Figure 4: Percentages of emotional child speech samples perceived by the listeners: correctly recognized age with the perception rate of 0-0.24 (gray), with the rate of 0.25-0.49 (white color), with the rate of 0.5-0.74 (sloping hatch) and with the rate of 0.75-1.0 (black)

297 tral state, and comfort state, respectively. The listeners recognized child's age
 298 mainly with a perception rate of 0.5-0.74 (Figure 4). Generally the agreement
 299 among the listeners in determining child's age was less than in gender and
 300 emotional state recognition. The amount of neutral samples on which most of
 301 listeners (0.75-1.0) correctly recognized child's age was higher than comfort and
 302 discomfort samples. This fact can be explained from the point of view of the
 303 pitch values. Decrease of pitch values from 3 to 7 years is shown. We can assume
 304 that listeners first of all lean on the pitch values in the determining the age of
 305 the child. Pitch values in emotional speech higher than in neutral speech that
 306 can confuse listeners.

307 The main sources of error in the age recognition task were the following:
 308 the listeners confused the 5 year old group with other ages (close ages were
 309 confused); speech samples of 5 year old children uttered in the discomfort and
 310 comfort states were often attributed to 3-4 year old kids. Confusion matrices
 311 are presented in Table 6. Based on these results, we can conclude that the age
 312 prediction task is more difficult for humans than the gender prediction task in
 313 our database.

314 Recently, two studies on human and machine prediction of gender and age of
 315 children are presented by (Safavi et al., 2013, 2014). These studies use material

Table 6: Confusion matrices for child’s age prediction in three emotional states.

	Emotional state								
	discomfort			neutral			comfort		
Age	3-4 y	5 y	6-7 y	3-4 y	5 y	6-7 y	3-4 y	5 y	6-7 y
3-4 y	53	39	8	65	29	6	69	28	3
5 y	35	49	16	22	50	28	31	44	25
6-7 y	16	36	48	13	44	43	17	43	40

316 from an age range of 5 to 13 years, which does not entirely overlap with the
 317 child age range in our data. Safavi et al. (2013) analyze the spectrum (24
 318 filtered frequencies) obtained from sliding windows of length 20ms shifted with
 319 10 ms. It was shown that the frequencies below 1.8 kHz and above 3.8 kHz are
 320 most useful for gender identification for older children (13-16 years), and the
 321 frequencies above 1.4 kHz are most useful for the youngest children (5-9 years).
 322 The frequencies above 5.5 kHz are the least useful for age identification (Safavi
 323 et al., 2014). Results show that in both cases, machine identification is more
 324 accurate compared to humans.

325 In the following machine learning experiments, we use a standard set of
 326 supra-segmental features without focus on feature or frequency band selection.
 327 Finding the most potent spectrum/cepstrum bands and most predictive features
 328 for these three classification tasks are left as future work.

329 4. Automatic Classification Systems for Paralinguistic Analysis

330 In this section, we investigate machine classification of the emotion, age,
 331 and gender of the child from speech segments. While there are several studies
 332 for the automatic processing of child speech (e.g. (Potamianos et al., 2011;
 333 Meinedo and Trancoso, 2011; Bolaños et al., 2011; Safavi et al., 2013, 2014),
 334 etc.), automatic detection and classification of emotional states of speech of
 335 children in natural conditions is a new direction of research Batliner et al. (2008);
 336 Lyakso et al. (2015). As our previous perception analysis reveals, the recognition

337 of children’s emotions from speech is hard, even though some prosodic patterns
338 can be discerned. The overall human recognition accuracy for a balanced three-
339 class problem (i.e. discomfort, comfort, neutral) is found as 65%.

340 In the following subsections, we provide a brief overview of paralinguistic
341 analysis and major elements of its pipeline, followed by the experimental results.

342 *4.1. Background on Automatic Paralinguistic Speech Analysis*

343 Paralinguistics (meaning alongside linguistics), studies short term states and
344 long term traits of the speaker(s), particularly focusing on non-verbal aspects
345 of speech that convey emotions. It deals with how the words are spoken, rather
346 than what is being spoken. Speech Emotion Recognition (SER) is a branch of
347 paralinguistics, related to speaker state and trait recognition (SSTR).

348 A general pipeline for SSTR is shown in Figure 5 (Schuller, 2011). The speech
349 signal is processed to extract informative features, which are fed to machine
350 learning modules for acoustic- and/or language-based classification. In paralin-
351 guistic analysis, acoustic models refer to affect classification models trained on
352 features derived from acoustic/prosodic Low Level Descriptors (LLD) such as
353 pitch, energy, jitter, shimmer and MFCCs (Schuller, 2011). On the other hand,
354 language-based classification employs linguistic information provided by auto-
355 matic speech recognition (ASR), which for instance can be represented as a
356 bag-of-words. While emotion recognition is a worthy goal on its own, it can also
357 help ASR through emotion-specific models. In turn, emotion recognition per-
358 formance can be improved by good ASR, by providing robust linguistic features
359 to be fused with acoustic features (Schuller et al., 2004). However, this is not al-
360 ways possible, since the recognition of affective (emotional) speech is itself very
361 challenging (Schuller, 2011). In the present study, we use only acoustic models
362 due to lack of Russian ASR trained on child speech, and since ASR trained on
363 adult speech does not work on children’s speech. Working on Russian child ASR
364 is one of our future research directions.

365 The state-of-the-art computational paralinguistics systems use large scale
366 suprasegmental feature extraction via passing a set of summarizing statistical

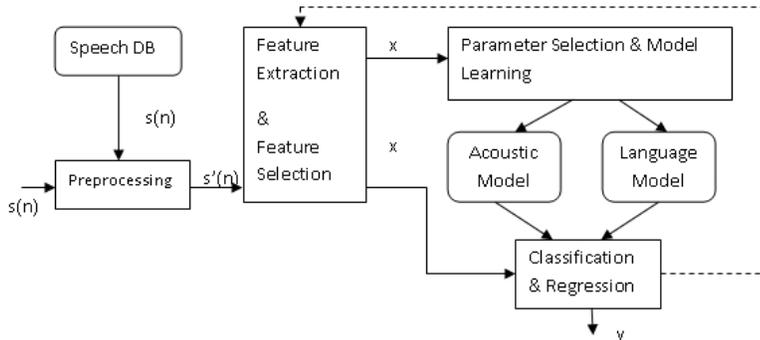


Figure 5: General speaker state and trait recognition pipeline.

367 functionals (such as moments, extremes) over LLD contours (Eyben et al., 2010).
 368 Pitch, Formants (resonant frequencies of vocal tract filter), Mel Frequency Cep-
 369 stral Coefficients (MFCC), Modulation Spectrum, Relative Spectral Transform
 370 - Perceptual Linear Prediction (RASTA-PLP), Energy and variation features
 371 (i.e. Shimmer and Jitter) are frequently used as LLDs (Schuller, 2011). Among
 372 these, MFCC and RASTA-PLP are the most widely used. In line with the state-
 373 of-the-art, we extract acoustic features in this work using the freely available
 374 openSMILE tool (Eyben et al., 2010), using a standard configuration file.

375 The most commonly employed classifiers in paralinguistics are Support Vec-
 376 tor Machines (SVM), Artificial Neural Networks (ANN), Gaussian Mixture
 377 Models (GMM), and Hidden Markov Models (HMM). The state-of-the-art mod-
 378 els of SER for the current databases are those trained with Support Vector Ma-
 379 chines (SVMs) and Deep Neural Networks (DNN) (Schuller, 2011). From the
 380 ANN family, Extreme Learning Machines (ELM), which combine fast model
 381 learning with accurate prediction capability, are recently applied to multi-modal
 382 emotion recognition and computational paralinguistics, obtaining state-of-the-
 383 art results with modest computational resources (Kaya et al., 2015a,b; Kaya and
 384 Salah, 2016). Consequently, we employ Kernel ELMs (Huang et al., 2012) in
 385 this work, as well as a fast and robust classifier based on Partial Least Squares
 386 (PLS) regression (Wold, 1985). In the next subsection, we give a brief summary
 387 of these two classification approaches. As a further baseline, we use SVMs.

388 *4.2. Background on Least Squares Regression based Classifiers*

389 To learn a classification model, we employ kernel extreme learning machine
 390 (ELM) and Partial Least Squares (PLS) regression due to their fast and accu-
 391 rate learning capability (Wold, 1985; Huang et al., 2012) and state-of-the-art
 392 achievements in recent audio and video based challenges (Kaya et al., 2015b;
 393 Kaya and Karpov, 2016; Kaya et al., 2017).

ELM proposes a Single Layer Feedforward Network (SLFN) architecture, but unsupervised, even random generation of the hidden node output matrix $\mathbf{H} \in \mathbb{R}^{N \times h}$, where N and h denote the number of data samples and the hidden neurons, respectively. The hidden node output matrix \mathbf{H} is obtained by projecting the data matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$ using a randomly generated first layer weight matrix $\mathbf{W} \in \mathbb{R}^{d \times h}$, and subsequently applying an infinitely differentiable non-linear activation function (e.g. logistic sigmoid). The actual learning takes place in the second layer between \mathbf{H} and the label matrix $\mathbf{T} \in \mathbb{R}^{N \times L}$, where L is the number of classes. Let y^t denote the class label of t^{th} data instance, in the case of L -class classification, \mathbf{T} is represented in one vs. all coding as follows:

$$\mathbf{T}_{t,l} = \begin{cases} +1 & \text{if } y^t = l, \\ -1 & \text{if } y^t \neq l. \end{cases} \quad (1)$$

394 In case of regression, $\mathbf{T} \in \mathbb{R}^{N \times 1}$ is the continuous target variable.

The second level weights $\beta \in \mathbb{R}^{h \times L}$ are learned by least squares solution to a set of linear equations $\mathbf{H}\beta = \mathbf{T}$. The output weights can be learned via:

$$\beta = \mathbf{H}^\dagger \mathbf{T}, \quad (2)$$

where \mathbf{H}^\dagger is the Moore-Penrose generalized inverse (Rao and Mitra, 1971) that gives the minimum L_2 norm solution to $\|\mathbf{H}\beta - \mathbf{T}\|$, simultaneously minimizing the norm of $\|\beta\|$. To increase the robustness and generalization capability, the optimization problem of ELM is reformulated using a regularization coefficient on the residual error $\|\mathbf{H}\beta - \mathbf{T}\|$. The learning rule of this alternative ELM is related to Least Square SVMs (LSSVM) via the following output weight learning formulation:

$$\beta = \mathbf{H}^T \left(\frac{\mathbf{I}}{C} + \mathbf{H}\mathbf{H}^T \right)^{-1} \mathbf{T}, \quad (3)$$

where \mathbf{I} is the $N \times N$ identity matrix, and C , which is used to regularize the linear kernel $\mathbf{H}\mathbf{H}^T$, corresponds to the complexity parameter of LSSVM (Suykens and Vandewalle, 1999). This formulation is further simplified by noting that the hidden layer matrix need not be generated explicitly given a kernel \mathbf{K} , which can be seen identical to Kernel Regularized Least Squares (Huang et al., 2012; Rifkin et al., 2003):

$$\beta = \left(\frac{\mathbf{I}}{C} + \mathbf{K}\right)^{-1}\mathbf{T}. \quad (4)$$

395 The second approach we use for classification is partial least squares (PLS)
 396 regression. PLS regression between two sets of variables $\mathbf{X} \in \mathbb{R}^{N \times d}$ and $\mathbf{Y} \in$
 397 $\mathbb{R}^{N \times p}$ is based on decomposing the matrices as $\mathbf{X} = \mathbf{U}_x\mathbf{V}_x + r_x$, $\mathbf{Y} = \mathbf{U}_y\mathbf{V}_y + r_y$,
 398 where \mathbf{U} denotes the latent factors, \mathbf{V} denotes the loadings and r stands for the
 399 residuals. The decomposition is done by finding projection weights $\mathbf{W}_x, \mathbf{W}_y$
 400 that jointly maximize the covariance of corresponding columns of $\mathbf{U}_x = \mathbf{X}\mathbf{W}_x$
 401 and $\mathbf{U}_y = \mathbf{Y}\mathbf{W}_y$. For further details of PLS regression, the reader is referred
 402 to (Wold, 1985). When PLS is applied to the classification problem in a one-
 403 versus-all setting, it learns the regression function between the feature matrix
 404 \mathbf{X} and the binary label vector \mathbf{Y} , and the class giving the highest regression
 405 score is taken as prediction. The number of latent factors is a hyper-parameter
 406 to tune via cross-validation.

407 4.3. Features and Performance Measures used in Machine Classification

408 We extract openSMILE (Eyben et al., 2010) features with a configuration
 409 file used in the INTERSPEECH 2010 Computational Paralinguistics Challenge
 410 (ComParE) as baseline set (Schuller et al., 2010a). This feature set contains
 411 1.582 suprasegmental features obtained by passing 21 descriptive functionals
 412 (e.g. moments, percentiles, regression coefficients) on 38 Low Level Descriptors
 413 (LLD) extracted from the speech signal (see Table 7 for details). This config-
 414 uration file is preferred over the one used in the 2015 edition of the ComParE
 415 Challenge, since in our recent work on a subset of this corpus (Lyakso et al.,
 416 2015), ComParE 2010 baseline set gave better results compared to the 2015
 417 version, which is a 6.373-dimensional acoustic feature set.

Table 7: The openSMILE feature set used in the study with a standard configuration from INTERSPEECH 2010 ComParE Challenge (Schuller et al., 2010a). DDP: difference of difference of periods; LSP: line spectral pairs.

Descriptors	Functionals
F0 by Sub-Harmonic Sum.	arithmetic mean, standard deviation
F0 Envelope	linear regression coefficients 1/2
Jitter DDP	linear regression error (quadratic/absolute)
Jitter local	percentile 1/99
Log Mel Freq. Band [0-7]	percentile range 99-1
LSP Frequency [0-7]	quartile 1/2/3
MFCC [0-14]	quartile range 2-1/3-2/3-1
PCM Loudness	relative position of minimum/maximum
Probability of Voicing	skewness, kurtosis
Shimmer local	up-level time 75/90

418 In all classification experiments reported below, the acoustic features are
419 first normalized and kernelized using Linear and Radial Basis Function (RBF)
420 kernels before classification. Kernelization refers to obtaining an instance sim-
421 ilarity matrix dubbed *kernel* $\mathbf{K} \in \mathbb{R}^{N \times N}$ from the data matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$. It
422 is popularly employed in linear classifiers to avoid the curse of dimensionality
423 (especially when $d \gg N$) and to allow non-linear separability in an implic-
424 itly mapped hyper-space by means of non-linear kernels, such as RBF. As pre-
425 processing, we apply z-normalization (i.e. standardization to zero-mean, unit
426 variance) or min-max normalization to $[0,1]$ range. The hyper parameters of
427 classifiers are optimized using a two-fold, speaker independent cross-validation
428 within the training set. The optimal parameters are finally used for model
429 training and predicting the labels of the test set, which is only used once for
430 reporting the results.

431 We report classification results in terms of accuracy and Unweighted Av-
432 erage Recall (UAR), which is introduced as performance measure in the IN-
433 TERSPEECH 2009 Emotion Challenge (Schuller et al., 2009). UAR is used to

434 overcome the biased calculation of accuracy towards the majority class. It also
 435 gives a chance-level baseline performance as $1/K$, where K is the number of
 436 classes. Therefore, in a 3-class problem, we have 33.3% chance-level UAR.

437 *4.4. Automatic Classification of Child Emotional States*

438 In affective computing, *arousal* and *valence* are the two main dimensions
 439 along which continuous and dimensional affect is measured (Russell, 1980).
 440 Arousal is defined as physiological/psychological state of being (re-)active, while
 441 valence is the feeling of positiveness (Schuller, 2011). The comfort classification
 442 can be thought as a three-state valence classification problem. It is well known
 443 that valence classification from acoustics is poorer compared to arousal classi-
 444 fication, and is almost at chance level in challenging conditions (i.e. without
 445 adaptation to cross-corpus settings) (Schuller et al., 2010b). The test set classi-
 446 fication results (in terms of both accuracy and UAR) obtained from two normal-
 447 ization, two kernel and three classifier alternatives are presented in Table 8. The
 448 confusion matrices for the predictions with the highest UAR (z-norm, Linear
 449 kernel, ELM) are shown in Table 9.

Table 8: Test set automatic classification results (%) for three emotional states.

Preprocess		UAR			Accuracy		
Normalization	Kernel	PLS	ELM	SVM	PLS	ELM	SVM
z-norm	Linear	47.9	51.5	48.6	52.0	52.0	53.7
min-max	Linear	47.2	49.1	48.2	51.8	53.9	53.3
z-norm	RBF	50.7	50.9	49.6	56.0	56.2	55.6
min-max	RBF	48.1	48.9	49.6	52.0	53.9	54.9

450 From Table 8, we observe that using only the acoustic features, it is possible
 451 to get higher than chance-level UAR scores. However, the best accuracy (56%)
 452 is much below the gold standard obtained from human perception experiments
 453 (66%). Note that the human perception experiments described in Section 3
 454 were done on a subset of the test set with 90 instances, and the decisions of 300

455 human listeners were fused. Here, we report the performance of individual clas-
 456 sifiers over 475 test set instances. Moreover, the classifiers do not benefit from
 457 linguistic information that might have been useful for human discrimination.

458 The results in Table 9 are quite different than their human perception coun-
 459 terpart. Interestingly, the automatic emotion recognition performances are
 460 found as 55.3% and 45.5% UAR for 3-4 years group and 6-7 years group, re-
 461 spectively. We observe that while the UAR performance of human perception
 462 of affective states improves with increasing age, acoustics-based automatic clas-
 463 sification gives higher performance with younger children. This may imply that
 464 the human listeners make implicit use of linguistic information that develops
 465 with age, such as the expanding vocabulary of the child.

Table 9: Row normalized (%) confusion matrices for automatic emotional state recognition giving the highest UAR (51.5% using z-norm, linear kernel, ELM) with age-group breakdown.

	Child’s age								
	3-4 years			5 years			6-7 years		
State	disc	neut	comf	disc	neut	comf	disc	neut	comf
discomfort	51	15	34	64	24	12	27	59	14
neutral	14	69	17	34	48	18	17	58	25
comfort	34	20	46	41	10	48	19	30	51

466 *4.5. Automatic Classification of Child Age Group and Gender*

467 The results for automatic three-level age group classification and the con-
 468 fusion matrices corresponding to the best UAR performance (z-norm, Linear
 469 kernel, PLS) are given in Tables 10 and 11, respectively. Comparing the UAR
 470 performances against the chance level, both in human and in machine classifi-
 471 cation this task is found to be the hardest among the three paralinguistic tasks
 472 dealt with in this paper. On the other hand, this is the only task where the best
 473 overall UAR score outperforms the one found in the human perception tests.

474 Considering the confusion matrices in Table 11, we observe that similar to
 475 the case with humans, confusions arise between the middle group (5 years) and

Table 10: Test set automatic classification results (%) for three-class age group estimation.

Preprocess		UAR			Accuracy		
Normalization	Kernel	PLS	ELM	SVM	PLS	ELM	SVM
z-norm	Linear	54.2	53.0	52.5	51.8	48.2	53.3
min-max	Linear	49.1	48.7	52.4	46.7	47.2	53.7
z-norm	RBF	52.0	48.9	52.4	48.0	48.0	49.9
min-max	RBF	48.4	50.0	51.5	46.1	48.2	48.6

Table 11: Confusion matrices for child age classification giving the highest test set UAR (54.2% using z-norm, Linear kernel and PLS) with comfort state breakdown.

	Emotional state								
	discomfort			neutral			comfort		
Age	3-4 y	5 y	6-7 y	3-4 y	5 y	6-7 y	3-4 y	5 y	6-7 y
3-4 y	64	25	11	48	45	7	63	31	6
5 y	12	72	16	8	64	28	12	57	31
6-7 y	0	73	27	8	54	38	13	41	46

476 the other two groups, which can be attributed to the narrow age span of this
477 class. Confusion between 3-4 and 6-7 years is low in general. On the other hand,
478 it is interesting to see that higher UAR performance of age prediction is observed
479 with comfort and discomfort states (55%) compared to the neutral state (50%).
480 Ordinarily, one would expect the contrary, as both in speech recognition and in
481 recognition of speaker traits, emotional speech generally gives lower performance
482 compared to neutral speech Schuller (2011).

483 Finally, the automatic classification results for child gender and the confusion
484 matrices corresponding to the best UAR performance prediction are listed in
485 Tables 12 and 13, respectively. For gender classification, the best machine UAR
486 performance is found lower than the UAR obtained from the human perception
487 test. Comparing the confusion matrices of machine and human classification, we
488 see a similar pattern in the discomfort state: recall of the male class is markedly
489 higher than the recall of the female class. On the other hand, the recall patterns
490 are different in the other two emotional states: the recall of female class is
491 higher than the recall of male class in machine classification, while it is the
492 opposite case for human perception. In the special case of discomfort, females
493 are highly confused as male (51% by machine, 48% by humans). It is widely
494 known that children’s acoustic features are very similar in two genders and
495 thus discrimination is hard. Our experimental results indicate that discomfort
496 vocalizations of female children resemble those of male children.

497 Mann-Whitney U tests are administered on predictions to measure statis-
498 tical significance of gender recall performance under different emotional states.
499 The test results indicate that female recall performances are statistically dif-
500 ferent between emotional states ($p < 10^{-6}$, $p < 10^{-4}$, $p < 0.05$ for neutral vs.
501 discomfort, neutral vs. comfort and comfort vs. discomfort, respectively). Male
502 recall is found significantly different between comfort and discomfort as well as
503 between neutral and discomfort states ($p < 0.05$), while no statistical difference
504 is observed between male recall performances under neutral and comfort states.

505

Table 12: Test set automatic classification results (%) for child’s gender.

Preprocess		UAR			Accuracy		
Normalization	Kernel	PLS	ELM	SVM	PLS	ELM	SVM
z-norm	Linear	54.6	49.5	52.4	56.6	49.1	52.6
min-max	Linear	57.0	50.2	51.7	58.5	50.5	52.0
z-norm	RBF	48.7	49.5	48.6	48.6	49.7	48.2
min-max	RBF	56.7	52.5	50.6	58.5	53.3	50.5

Table 13: Confusion matrices for automatic gender classification giving the highest UAR (57.0% using min-max normalization, Linear kernel and PLS) with emotional state breakdown.

	Emotional state					
	discomfort		neutral		comfort	
Gender	male	female	male	female	male	female
male	60	40	37	63	44	56
female	51	49	13	87	35	65

506 **5. Discussion and Conclusions**

507 The present work is part of an emotional development study, which inves-
508 tigates emotional states in verbal and non-verbal behavior of kids during the
509 first seven years of life. Choosing the age range as 3-7 years is due to the evo-
510 lution of the grammatical skills of speech at 4 years, and the ability of effective
511 communication of a child with an adult. In this age range, regulation of emo-
512 tional expressions is not fully developed yet, and the emotional expressions are
513 purer, as the contribution of society in the organization of the child’s behavior
514 is comparatively small. The upper bound age of seven years is associated with
515 the end of the preschool time of children in the Russian Federation. There are
516 only a few databases with emotional speech of children before 4 years of age in
517 general (Lyakso et al., 2010a).

518 The presented experimental results show that lexical information has more
519 discriminative power in recognition of comfort compared to acoustic features for

520 speech samples of 6-7 year old children. The database contains instances (es-
521 pecially for smaller infants) that are difficult to annotate, even by the parents.
522 Older infants speak more clearly. Human perception of emotion is higher in the
523 speech of older children, which has several implications. It is harder to recognize
524 sentiment with younger children, since linguistic and acoustic control skills are
525 not mature enough. Furthermore, the comparative experiments in age and gen-
526 der classification tasks also reveal the importance of linguistic models for better
527 discrimination. There is no ASR solution for Russian child speech yet, which
528 suggests this would be good future research direction, both as an independent
529 study, and for multi-modal paralinguistic analysis. Analyzing/monitoring child
530 emotion in early ages is important not only for linguistics, but also for anal-
531 ysis of neurological development/disorders. Moreover, all three paralinguistic
532 tasks targeted here are important for developing intelligent tutoring systems for
533 pre-school education.

534 In age group classification, the automatic classification study reveals higher
535 performance compared to human perception, which can be taken as a gold
536 standard, and falls below human performance in the other two. Especially
537 for emotional state classification, the human listeners may be using linguistic
538 content to their advantage, and it is difficult to quantify this. Further research
539 is necessary to achieve and to outperform human performance. Multimodal
540 fusion from linguistic and visual cues seems a promising step in this direction.
541 Another important research direction is cross-corpus, cross-language analysis in
542 child paralinguistics, as it is the key to leveraging additional data sources in
543 training automatic systems.

544 **6. Acknowledgments**

545 The work was supported by the Russian Foundation for Basic Research
546 (grants № 16-06-00024, 15-06-07852, and 16-37-60100), Russian Foundation for
547 Basic Research - DHSS (grant № 17-06-00503), by the grant of the President
548 of Russia (project № MD-254.2017.8), by the Government of Russia (grant №

549 074-U01), by Boğaziçi University (project BAP 16A01P4) and by the BAGEP
550 Award of the Science Academy.

551 **References**

552 Batliner, A., Blomberg, M., D’Arcy, S., Elenius, D., Giuliani, D., Gerosa, M.,
553 Hacker, C., Russell, M.J., Steidl, S., Wong, M., 2005. The PF_STAR chil-
554 dren’s speech corpus, in: Proc. INTERSPEECH, pp. 2761–2764.

555 Batliner, A., Steidl, S., Nöth, E., 2008. Releasing a thoroughly annotated and
556 processed spontaneous emotional database: the FAU Aibo Emotion Corpus,
557 in: Proc. LREC-2008 Workshop of on Corpora for Research on Emotion and
558 Affect, pp. 28–31.

559 Bell, L., Boye, J., Gustafson, J., Heldner, M., Lindström, A., Wirén, M., 2005.
560 The Swedish NICE Corpus–spoken dialogues between children and embodied
561 characters in a computer game scenario, in: Proc. EUROSPEECH, ISCA. pp.
562 2765–2768.

563 Bolaños, D., Cole, R.A., Ward, W., Borts, E., Svirsky, E., 2011. FLORA: Fluent
564 oral reading assessment of children’s speech. ACM Transactions on Speech
565 and Language Processing (TSLP) 7, 16.

566 Burkhardt, F., Paeschke, A., Rolfes, M., Sendlmeier, W.F., Weiss, B., 2005. A
567 database of German emotional speech, in: Proc. INTERSPEECH, pp. 1517–
568 1520.

569 Engberg, I.S., Hansen, A.V., 1996. Documentation of the Danish Emotional
570 Speech Database DES. Internal AAU report, Center for Person Kommunika-
571 tion, Denmark , 1–22.

572 Eyben, F., Wöllmer, M., Schuller, B., 2010. OpenSMILE: the Munich versatile
573 and fast open-source audio feature extractor, in: Proc. 18th ACM Int. Conf.
574 on Multimedia, ACM. pp. 1459–1462.

575 Fringi, E., Lehman, J.F., Russell, M., 2015. Evidence of phonological processes
576 in automatic recognition of children’s speech, in: Proc. INTERSPEECH, pp.
577 1621–1624.

578 Gerosa, M., Giuliani, D., Brugnara, F., 2007. Acoustic variability and automatic
579 recognition of children’s speech. *Speech Communication* 49, 847–860.

580 Govender, A., Wet, F.d., Tapamo, J.R., 2015. HMM adaptation for child speech
581 synthesis, in: Proc. INTERSPEECH, pp. 1640–1644.

582 Guo, J., Paturi, R., Yeung, G., Lulich, S.M., Arsikere, H., Alwan, A., 2015. Age-
583 dependent height estimation and speaker normalization for children’s speech
584 using the first three subglottal resonances, in: Proc. INTERSPEECH, pp.
585 1665–1669.

586 Huang, G.B., Zhou, H., Ding, X., Zhang, R., 2012. Extreme learning machine
587 for regression and multiclass classification. *IEEE Transactions on Systems,
588 Man, and Cybernetics, Part B: Cybernetics* 42, 513–529.

589 Jürgens, R., Drolet, M., Pirow, R., Scheiner, E., Fischer, J., 2013. Encoding
590 conditions affect recognition of vocally expressed emotions across cultures.
591 *Frontiers in Psychology* 4, 111.

592 Kaya, H., Gürpınar, F., Afshar, S., Salah, A.A., 2015a. Contrasting and com-
593 bining least squares based learners for emotion recognition in the wild, in:
594 Proceedings of the 2015 ACM on International Conference on Multimodal
595 Interaction, ACM. pp. 459–466.

596 Kaya, H., Gürpınar, F., Salah, A.A., 2017. Video-based emotion recognition
597 in the wild using deep transfer learning and score fusion. *Image and Vision
598 Computing* doi:<http://dx.doi.org/10.1016/j.imavis.2017.01.012>.

599 Kaya, H., Karpov, A.A., 2016. Fusing acoustic feature representations for com-
600 putational paralinguistics tasks, in: Proc. INTERSPEECH, San Francisco,
601 USA. pp. 2046–2050.

- 602 Kaya, H., Karpov, A.A., Salah, A.A., 2015b. Fisher vectors with cascaded
603 normalization for paralinguistic analysis, in: Proc. INTERSPEECH, Dresden,
604 Germany. pp. 909–913.
- 605 Kaya, H., Salah, A.A., 2016. Combining modality-specific extreme learning ma-
606 chines for emotion recognition in the wild. *Journal on Multimodal User Inter-*
607 *faces* 10, 139–149. doi:<http://dx.doi.org/10.1007/s12193-015-0175-6>.
- 608 Kaya, H., Salah, A.A., Gurgen, S.F., Ekenel, H., 2014. Protocol and baseline
609 for experiments on Bogazici University Turkish emotional speech corpus, in:
610 22nd IEEE Signal Processing and Communications Applications Conference
611 (SIU), pp. 1698–1701.
- 612 Laukka, P., Elfenbein, H.A., Söder, N., Nordström, H., Althoff, J., Chui, W.,
613 Iraki, F.K., Rockstuhl, T., Thingujam, N.S., 2013. Cross-cultural decoding
614 of positive and negative non-linguistic emotion vocalizations. *Frontiers in*
615 *Psychology* 4, 353.
- 616 Liao, H., Pundak, G., Siohan, O., Carroll, M.K., Coccaro, N., Jiang, Q.M.,
617 Sainath, T.N., Senior, A., Beaufays, F., Bacchiani, M., 2015. Large vocab-
618 ulary automatic speech recognition for children, in: Proc. INTERSPEECH,
619 pp. 1611–1615.
- 620 Lindquist, K.A., Wager, T.D., Kober, H., Bliss-Moreau, E., Barrett, L.F., 2012.
621 The brain basis of emotion: a meta-analytic review. *Behavioral and Brain*
622 *Sciences* 35, 121–143.
- 623 Lyakso, E., Frolova, O., 2015. Emotion state manifestation in voice features:
624 Chimpanzees, human infants, children, adults, in: Proc. International Con-
625 ference on Speech and Computer (SPECOM). Springer, pp. 201–208.
- 626 Lyakso, E., Frolova, O., Dmitrieva, E., Grigorev, A., Kaya, H., Salah, A.A.,
627 Karpov, A., 2015. EmoChildRu: emotional child Russian speech corpus,
628 in: Proc. International Conference on Speech and Computer (SPECOM).
629 Springer, pp. 144–152.

- 630 Lyakso, E., Gromova, A., 2005. The acoustic characteristics of russian vowels
631 in children of 4-5 years of age. *Psychology of Language and Communication*
632 9, 5–14.
- 633 Lyakso, E., Kurazova, A., Gromova, A., Ostroukhov, A., 2006. Recognition of
634 words and phrases of 4-5-years-old children by adults, in: *Proc. International*
635 *Conference on Speech and Computer (SPECOM)*, pp. 567–570.
- 636 Lyakso, E.E., Frolova, O.V., Grigorev, A.S., Sokolova, V.D., Yarotskaya, K.A.,
637 2016. Recognition of adults emotional state of typically developing children
638 and children with autism spectrum disorders. *Neuroscience and Behavioral*
639 *Physiology* 102, 729–741.
- 640 Lyakso, E.E., Frolova, O.V., Kurazhova, A.V., Gaikova, J.S., 2010a. Russian
641 infants and children’s sounds and speech corpuses for language acquisition
642 studies, in: *Proc. INTERSPEECH*, pp. 1981–1888.
- 643 Lyakso, E.E., Ushakova, T.N., Frolova, O.V., Kurazhova, A.V., Bednaya, E.D.,
644 Gaikova, J.S., Grigoriev, A.S., Soloviev, A.N., Ostrouchov, A.V., 2010b. Rus-
645 sian children’s vocabulary, speech imitation and reading skills mastery. *Inter-*
646 *national Journal of Psychophysiology* 77, 310.
- 647 Meinedo, H., Trancoso, I., 2011. Age and gender detection in the I-DASH
648 project. *ACM Transactions on Speech and Language Processing (TSLP)* 7,
649 13.
- 650 Pérez-Espinosa, H., Reyes-García, C.A., Villaseñor-Pineda, L., 2011. EmoWis-
651 consin: an emotional children speech database in Mexican Spanish, in: *Proc.*
652 *Affective Computing and Intelligent Interaction*, LNCS vol. 6975. Springer,
653 pp. 62–71.
- 654 Potamianos, A., Giuliani, D., Narayanan, S.S., Berkling, K., 2011. Introduction
655 to the special issue on speech and language processing of children’s speech
656 for child-machine interaction applications. *ACM Transactions on Speech and*
657 *Language Processing (TSLP)* 7, 11.

- 658 Rao, C.R., Mitra, S.K., 1971. Generalized inverse of matrices and its applica-
659 tions. volume 7. Wiley New York.
- 660 Rifkin, R., Yeo, G., Poggio, T., 2003. Regularized least-squares classification.
661 Nato Science Series Sub Series III Computer and Systems Sciences 190, 131–
662 154.
- 663 Rigoulot, S., Wassiliwizky, E., Pell, M., 2013. Feeling backwards? how temporal
664 order in speech affects the time course of vocal emotion recognition. *Frontiers*
665 *in psychology* 4, 367–367.
- 666 Russell, J.A., 1980. A circumplex model of affect. *Journal of Personality and*
667 *Social Psychology* 39, 1161–1178.
- 668 Safavi, S., Jancovic, P., Russell, M.J., Carey, M.J., 2013. Identification of gender
669 from children’s speech by computers and humans., in: *INTERSPEECH*, Lyon,
670 France. pp. 2440–2444.
- 671 Safavi, S., Russell, M.J., Jancovic, P., 2014. Identification of age-group from
672 children’s speech by computers and humans., in: *INTERSPEECH*, Singapore.
673 pp. 243–247.
- 674 Schuller, B., 2011. Voice and speech analysis in search of states and traits, in:
675 *Computer Analysis of Human Behavior*. Springer, pp. 227–253.
- 676 Schuller, B., Rigoll, G., Lang, M., 2004. Speech emotion recognition combin-
677 ing acoustic features and linguistic information in a hybrid support vector
678 machine-belief network architecture, in: *Proc. IEEE Int. Conf. on Acoustics,*
679 *Speech, and Signal Processing (ICASSP’04)*, IEEE, Montreal, Canada. pp.
680 577–580.
- 681 Schuller, B., Steidl, S., Batliner, A., 2009. The *INTERSPEECH 2009 Emotion*
682 *Challenge*, in: *Proc. INTERSPEECH*, pp. 312–315.
- 683 Schuller, B., Steidl, S., Batliner, A., Burkhardt, F., Devillers, L., Müller, C.A.,
684 Narayanan, S.S., 2010a. The *INTERSPEECH 2010 Paralinguistic Challenge*,
685 in: *Proc. INTERSPEECH*, pp. 2795–2798.

- 686 Schuller, B., Vlasenko, B., Eyben, F., Wöllmer, M., Stuhlsatz, A., Wendemuth,
687 A., Rigoll, G., 2010b. Cross-corpus acoustic emotion recognition: variances
688 and strategies. *IEEE Transactions on Affective Computing* 1, 119–131.
- 689 Suykens, J.A., Vandewalle, J., 1999. Least squares support vector machine
690 classifiers. *Neural processing letters* 9, 293–300.
- 691 Ververidis, D., Kotropoulos, C., 2006. Emotional speech recognition: Resources,
692 features, and methods. *Speech Communication* 48, 1162–1181.
- 693 Watson, R., Latinus, M., Noguchi, T., Garrod, O., Crabbe, F., Belin, P., 2014.
694 Crossmodal adaptation in right posterior superior temporal sulcus during
695 face–voice emotional integration. *The Journal of Neuroscience* 34, 6813–6821.
- 696 Wold, H., 1985. Partial least squares, in: Kotz, S., Johnson, N.L. (Eds.), *Ency-*
697 *clopedia of Statistical Sciences*. Wiley New York, pp. 581–591.
- 698 Yildirim, S., Narayanan, S., Potamianos, A., 2011. Detecting emotional state
699 of a child in a conversational computer game. *Computer Speech & Language*
700 25, 29–44.